**Pre-requisites:**

Last week's videos lectures and homework assignment.

**Recap:**

Thanks to Rayan Wong for summarizing it for us. ( https://medium.com/@taggatle/01-reinforcement-learning-move-37-introduction-c3449dac2d54 )

Now that we are familiar with the basics, it's time to move forward. Let us revisit the same problem we discussed last week. Robotic vacuum cleaner famously known as Roomba is a machine that cleans the floor. Roomba needs to start, clean, avoid obstacles and find the charging station. These 4 states describe the possible positions of the robot and the action describes the direction of motion. The robot can move left/right/up/down. The first (Battery Full) and the final (Charging) states are the terminal states. The goal is to find an optimal policy that maximizes the return from any initial states.
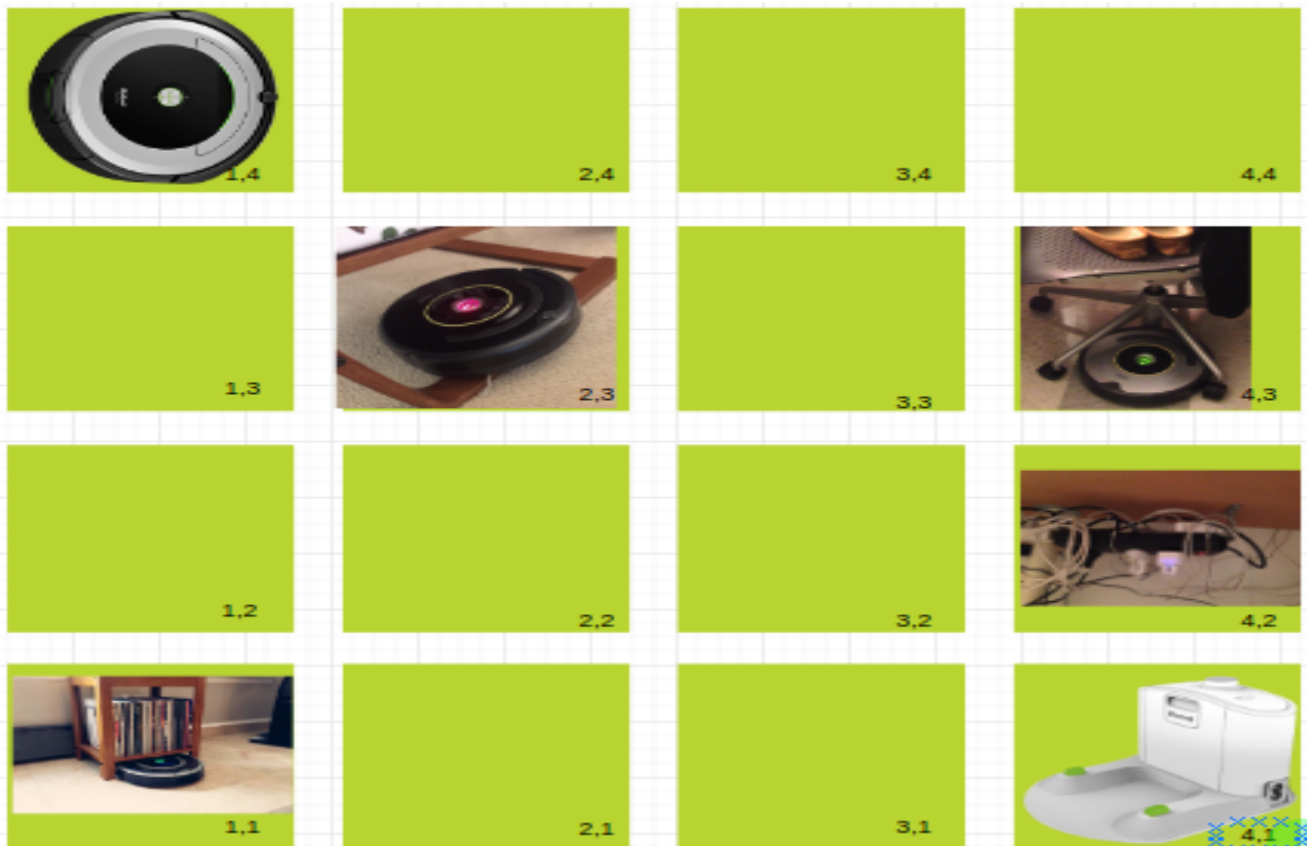


Fig 1

Above figure shows us the map of possible states that the Roomba can navigate through in form of a grid. For each step, Roomba gets a reward of 0 except when it reaches the goal to get a reward value of 1.

Let us use OpenAI gym and try to solve this problem. FrozenLake-v0 environment will be a best fit for this problem. The reward and next state models are stochastic p ($r\_t+1$ | $s\_t$ , $a\_t$ ), p($s\_t+1$ | $s\_t$ , $a\_t$ ) for this environment. To solve this problem, we should find an optimal policy that can reach the goal more than 70% of the times. (Please try to come up with a solution before looking at the solutions)

**Let's pull our sleeves and try a brute force method:**

We got 16 states and 4 possible moves that gives us 4^16 = 4294967296 possible policies to choose from. It is a computationally intensive task to evaluate all of them so we are going to choose few cases randomly and select the best among them.

```python
1.  import numpy
2.  import time
3.  import gym
4.  """
5.      Args:
6.          policy: [S, A] shaped matrix representing the policy.
7.          env: OpenAI gym env.
8.          render: boolean to turn rendering on/off.
9.  """
10. #Execution
11. def execute(env, policy, episodeLength=100, render=False):
12.     totalReward = 0
13.     start = env.reset()
14.     for t in range(episodeLength):
15.         if render:
16.             env.render()
17.         action = policy[start]
18.         start, reward, done, _ = env.step(action)
19.         totalReward += reward
20.         if done:
21.             break
22.     return totalReward
23.
24. #Evaluation
25. def evaluatePolicy(env, policy, n_episodes=100):
26.     totalReward = 0.0
27.     for _ in range(n_episodes):
28.         totalReward += execute(env, policy)
29.     return totalReward / n_episodes
30.
31. #Function for a random policy
32. def gen_random_policy():
33.     return numpy.random.choice(4, size=((16)))
34.
```

```
35. if __name__ == '__main__':
36.     env = gym.make('FrozenLake-v0')
37.     ## Policy search
38.     n_policies = 1000
39.     startTime = time.time()
40.     policy_set = [gen_random_policy() for _ in range(n_policies)]
41.     policy_score = [evaluatePolicy(env, p) for p in policy_set]
42.     endTime = time.time()
43.     print("Best score = %0.2f. Time taken = %4.4f seconds" %(numpy.max(policy_score) ,
        endTime - startTime))
```

**Output for the above script:**

**<u>Best score = 0.44. Time taken = 10.2633 seconds</u>**

The above script searches the environment for best policy in a random set of 1000 solutions and evaluates them. The best policy score we get is 0.44 in 10.2633 seconds. It means that the chance of agent reaching the goal is 44%. It is not even close to our goal. Random search does not work well for complex problems where the search space is huge.

How can we get this better???  Let's Dive Deeper...

The goal of the Roomba is to pick the best policy that will maximize the total rewards received from the environment.
$t_0$: Roomba observes the environment state $s_0$ and picks an action $a_0$, then upon performing its action, environment state becomes $s_1$ and the Roomba receives a reward $r_1$.

$t_1$: Roomba observes current state $s_1$ and picks an action $a_1$, then upon acting its action, environment state becomes $s_2$ and it receives a reward $r_2$.

$t_2$: Roomba observes current state $s_2$ and picks an action $a_2$, then upon acting its action, environment state becomes $s_3$ and it receives a reward $r_3$.

Thus the total reward received by the Roomba in response to the actions selected by its policy is going to be:

**Total reward = r_1 + r_2 + r_3 + r_4 + r_5 + ...**

However, it is common to use a discount factor to give higher weight if the rewards are nearby and lower weight if the rewards are far in the future.

**Total discounted reward = r_1 + $\gamma$ r_2 + $\gamma^2$ r_3 + $\gamma^3$ r_4 + $\gamma^4$ r_5+ …**

$$total\ discounted\ reward = \sum_{i=1}^{T} \gamma^{i-1} r_i$$

*T is the length of the episode. It can be infinity if there is huge length for the episode.*

Why discount factor???

The idea of using discount factor is to prevent the total reward from going to infinity (because $0 <= \gamma <= 1$). It also models the agent behavior when the agent prefers immediate rewards than rewards that are potentially received later in the future.

In our example Roomba can take two paths to reach goal state; one of them is longer but gives higher reward and there is also a shorter path with smaller reward. We need the Roomba to cover and clean maximum distance before it reaches the goal. By adjusting the $\gamma$ value we can control which path Roomba should prefer.

**Value Function:**

This term is very popular in reinforcement learning. It's denoted as **V(s)** and this represents how good a state is for an agent to be in. If an agent uses a given policy $\pi$ to select actions, then the corresponding value function is given by:

$$V^\pi(s) = \mathbb{E}[\sum_{i=1}^{T} \gamma^{i-1} r_i] \quad \forall s \in \mathbb{S}$$

*$\forall$ - Universal Quantification (https://en.wikipedia.org/wiki/Universal_quantification)*

Amongst all the possible value functions there will be an optimal value function that has higher value than other functions for all states.

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

The optimal policy **π\*** is the policy that corresponds to the optimal value function.

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

**Q Function:**

Q is a function of a state-action pair and returns real value. " **Q : S x A --> R** ". An optimal Q-function **Q\*(s,a)** means the expected total reward received by an agent starting in a state **s** and picks an action **a** will then move optimally forward. This can also be put as the indication for how good it is for an agent to pick action **a** while being in state **s.**

Since **V\*(s)** is the maximum expected total reward when starting from state a, it will be the maximum of Q\*(s,a) over all possible actions.

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

When a known optimal Q-function Q\*(s,a) exists, the optimal policy can be attained by choosing the action a that declares the maximum Q\*(s,a) for a particular state *s.*

$$\pi^*(s) = \arg\max_a Q^*(s, a) \quad \forall s \in \mathbb{S}$$

Bellman equation using dynamic programming gives a recursive definition for the optimal Q-function. The Q*(s,a) equals the summation of immediate reward after performing action *a* while in state *s* and the discounted expected future reward after transition to a next state *s'* .

$$Q^*(s, a) = R(s, a) + \gamma \, \mathbb{E}_{s'}[V^*(s')]$$
$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in \mathbb{S}} p(s'|s, a)V^*(s')$$

Since,
$$V^*(S) = \max_a Q^*(s, a)$$

$$V^*(S) = \max_a \left[ R(s, a) + \gamma \sum_{s' \in \mathbb{S}} p(s'|s, a)V^*(s') \right]$$

*\*\*please refer to last week's portion for more details.*

***Value Iteration:***

Value Iteration computes the optimal state value function by iteratively improving the estimate of **V(s)**. The algorithm initializes **V(s)** to arbitrary random values then it repeatedly updates the **Q(s,a)** and **V(s)** values until they converge. Value iteration is guaranteed to converge to the optimal values.

---

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathbb{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
| $\Delta \leftarrow 0$
| Loop for each $s \in \mathbb{S}$:
|     $v \leftarrow V(s)$
|     $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
|     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
    $\pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

---

Fig 2. *credits: Richard S. Sutton and Andrew G. Barto*

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})}\left[V^{\pi}(\mathbf{s}')\right]$$



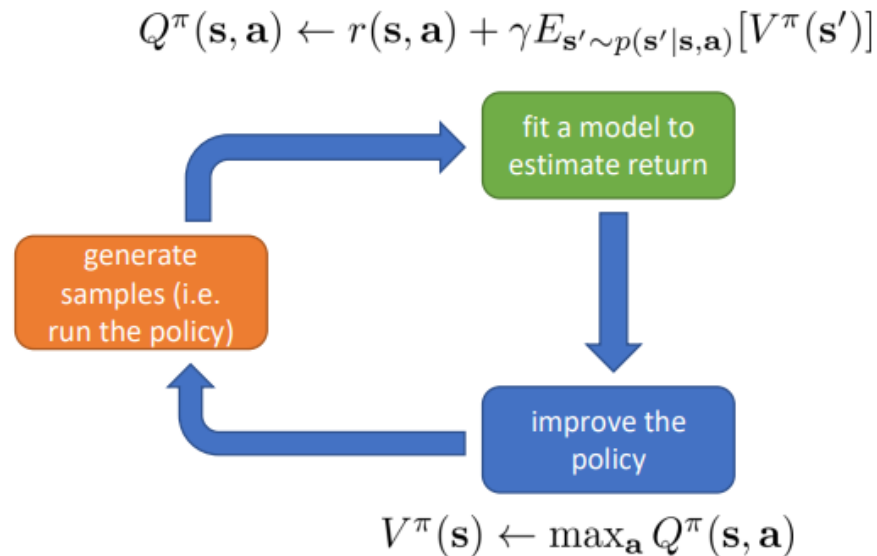$$V^{\pi}(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a})$$

*Fig 3. credits: Prof.Sergey Levine*

**Fun part now:** *let's take the above problem and try to solve it using  Value Iteration method.*

1. **import** numpy
2. **import** gym
3. **import** time
4. """
5.       Args:
6.             policy: [S, A] shaped matrix representing the policy.
7.             env: OpenAI gym env.
8.               env.P represents the transition probabilities of the environment.
9.               env.P[s][a] is a list of transition tuples (prob, next_state, reward, done).
10.             env.nS is a number of states in the environment.
11.             env.nA is a number of actions in the environment.
12.         gamma: Gamma discount factor.
13.         render: boolean to turn rendering on/off.
14. """
15. **def** execute(env, policy, gamma=1.0, render=False):
16.     start = env.reset()
17.     totalReward = 0
18.     stepIndex = 0
19.     **while** True:
20.         **if** render:
21.             env.render()

```python
22.        start, reward, done, _ = env.step(int(policy[start]))
23.        totalReward += (gamma ** stepIndex * reward)
24.        stepIndex += 1
25.        if done:
26.            break
27.    return totalReward
28.
29. # Evaluates a policy by running it n times.returns:average total reward
30. def evaluatePolicy(env, policy, gamma=1.0, n=100):
31.    scores = [
32.        execute(env, policy, gamma=gamma, render=False)
33.        for _ in range(n)]
34.    return numpy.mean(scores)
35.
36. # choosing the policy given a value-function
37. def calculatePolicy(v, gamma=1.0):
38.    policy = numpy.zeros(env.env.nS)
39.    for s in range(env.env.nS):
40.        q_sa = numpy.zeros(env.action_space.n)
41.        for a in range(env.action_space.n):
42.            for next_sr in env.env.P[s][a]:
43.                # next_sr is a tuple of (probability, next state, reward, done)
44.                p, s_, r, _ = next_sr
45.                q_sa[a] += (p * (r + gamma * v[s_]))
46.        policy[s] = numpy.argmax(q_sa)
47.    return policy
48.
49. # Value Iteration Agorithm
50. def valueIteration(env, gamma=1.0):
51.    value = numpy.zeros(env.env.nS)  # initialize value-function
52.    max_iterations = 10000
53.    eps = 1e-20
54.    for i in range(max_iterations):
55.        prev_v = numpy.copy(value)
56.        for s in range(env.env.nS):
```

```
57.        q_sa = [sum([p * (r + prev_v[s_]) for p, s_, r, _ in env.env.P[s][a]]) for a in
    range(env.env.nA)]
58.        value[s] = max(q_sa)
59.    if (numpy.sum(numpy.fabs(prev_v - value)) <= eps):
60.        print('Value-iteration converged at # %d.' % (i + 1))
61.        break
62.    return value
63.
64.if __name__ == '__main__':
65.    gamma = 1.0
66.    env = gym.make("FrozenLake-v0")
67.    optimalValue = valueIteration(env, gamma);
68.    startTime = time.time()
69.    policy = calculatePolicy(optimalValue, gamma)
70.    policy_score = evaluatePolicy(env, policy, gamma, n=1000)
71.    endTime = time.time()
72.    print("Best score = %0.2f. Time taken = %4.4f seconds" % (numpy.mean(policy_score),
    endTime - startTime))
```

**Output:**

Value-iteration converged at # 1373.

**Best score = 0.76. Time taken = 0.3675 seconds.**

**Policy Iteration:**

The value-iteration algorithm keeps improving the value function at each iteration, until the value-function converges. Since the agent only cares about the finding the optimal policy, sometimes the optimal policy will converge before the value function. Therefore, we have another algorithm called policy-iteration. Instead of repeatedly improving the value-function estimate, it will re-define the policy at each step and compute the value according to this new policy until the policy converges. Policy iteration is also guaranteed to converge to the optimal policy and it often takes less iterations to converge than the value-iteration algorithm.

Fig 4. *credits: Richard S. Sutton and Andrew G. Barto*



policy iteration:
1. evaluate $V^\pi(\mathbf{s})$
2. set $\pi \leftarrow \pi'$

$$\pi'(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg\max_{\mathbf{a}_t} A^\pi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

policy evaluation:

$$V^\pi(\mathbf{s}) \leftarrow r(\mathbf{s}, \pi(\mathbf{s})) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \pi(\mathbf{s}))}[V^\pi(\mathbf{s}')]$$

estimate $V^\pi$

fit a model to estimate return

generate samples (i.e. run the policy)

improve the policy

$\pi \leftarrow \pi'$

| | | | |
|---|---|---|---|
| 0.2 | 0.3 | 0.4 | 0.3 |
| 0.3 | 0.3 | 0.5 | 0.3 |
| 0.4 | 0.4 | 0.6 | 0.4 |
| 0.5 | 0.5 | 0.7 | 0.5 |

16 states, 4 actions per state

can store full $V^\pi(\mathbf{s})$ in a table!
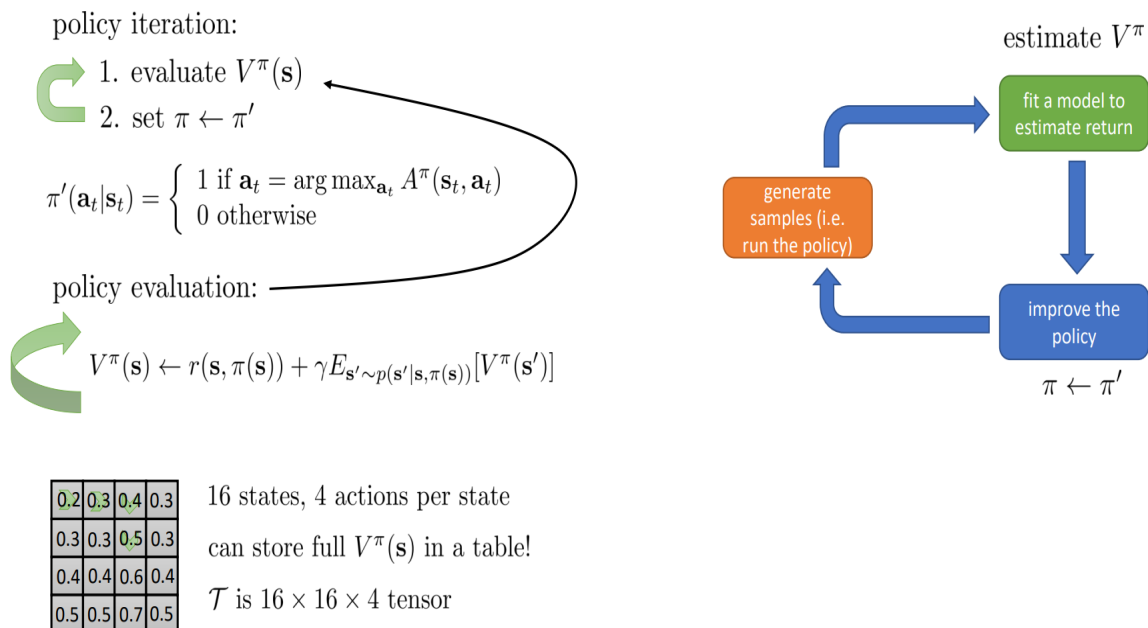
$\mathcal{T}$ is $16 \times 16 \times 4$ tensor

Fig 5. *credits: Prof.*Sergey Levine

**Fun part again:** *let's take the same problem and try to solve it using Policy Iteration method.*

73. **import** numpy

```python
import gym
import time
"""
    Args:
        policy: [S, A] shaped matrix representing the policy.
        env: OpenAI gym env.
            env.P represents the transition probabilities of the environment.
            env.P[s][a] is a list of transition tuples (prob, next_state, reward, done).
            env.nS is a number of states in the environment.
            env.nA is a number of actions in the environment.
        gamma: Gamma discount factor.
        render: boolean to turn rendering on/off.
"""
#executes an episode
def execute(env, policy, gamma = 1.0, render = False):
    start = env.reset()
    totalReward = 0
    stepIndex = 0
    while True:
        if render:
            env.render()
        start, reward, done , _ = env.step(int(policy[start]))
        totalReward += (gamma ** stepIndex * reward)
        stepIndex += 1
        if done:
            break
    return totalReward

def evaluatePolicy(env, policy, gamma = 1.0, n = 100):
    scores = [execute(env, policy, gamma, False) for _ in range(n)]
    return numpy.mean(scores)

#Extract the policy given a value-function
def extractPolicy(v, gamma = 1.0):
    policy = numpy.zeros(env.env.nS)
    for s in range(env.env.nS):
```

```python
110.        q_sa = numpy.zeros(env.env.nA)
111.        for a in range(env.env.nA):
112.            q_sa[a] = sum([p * (r + gamma * v[s_]) for p, s_, r, _ in  env.env.P[s][a]])
113.        policy[s] = numpy.argmax(q_sa)
114.    return policy
115.

116.#Iteratively calculates the value-function under policy.
117.def CalcPolicyValue(env, policy, gamma=1.0):
118.    value = numpy.zeros(env.env.nS)
119.    eps = 1e-10
120.    while True:
121.        previousValue = numpy.copy(value)
122.        for states in range(env.env.nS):
123.            policy_a = policy[states]
124.            value[states] = sum([p * (r + gamma * previousValue[s_]) for p, s_, r, _ in
     env.env.P[states][policy_a]])
125.        if (numpy.sum((numpy.fabs(previousValue - value))) <= eps):
126.            # value converged
127.            break
128.    return value
129.

130.#PolicyIteration algorithm
131.def policyIteration(env, gamma = 1.0):
132.    policy = numpy.random.choice(env.env.nA, size=(env.env.nS))  # initialize a random policy
133.    maxIterations = 1000
134.    gamma = 1.0
135.    for i in range(maxIterations):
136.        oldPolicyValue = CalcPolicyValue(env, policy, gamma)
137.        newPolicy = extractPolicy(oldPolicyValue, gamma)
138.        if (numpy.all(policy == newPolicy)):
139.            print ('Policy Iteration converged at %d.' %(i+1))
140.            break
141.        policy = newPolicy
142.    return policy
143.

144.if __name__ == '__main__':
```

```
145.    env_name  = 'FrozenLake-v0'
146.    env = gym.make(env_name)
147.    start = time.time()
148.    optimalPolicy = policyIteration(env, gamma = 1.0)
149.    scores = evaluatePolicy(env, optimalPolicy, gamma = 1.0)
150.    end = time.time()
151.    print("Best score = %0.2f. Time taken = %4.4f seconds" %(numpy.max(scores) , end -
        start))
```

**Output:**

**Policy-Iteration converged at step 3.**

**Best score = 0.78. Time taken = 0.1491 seconds.**

**Value-Iteration vs Policy-Iteration:**
Both value-iteration and policy-iteration algorithms can be used for *offline planning* where the agent is assumed to have prior knowledge about the effects of its actions on the environment (they assume the MDP model is known). Comparing each other, policy-iteration is computationally efficient as it often takes considerably fewer number of iterations to converge although each iteration is more computationally expensive.

Let us make it better next week!!!

**Credits:**

http://learning.mpi-sws.org/mlss2016/slides/2016-MLSS-RL.pdf

http://incompleteideas.net/book/bookdraft2018jan1.pdf

http://rail.eecs.berkeley.edu/deeprlcourse/static/slides/lec-7.pdf

https://gym.openai.com/evaluations/eval_4VyQBhXMRLmG9y9MQA5ePA

http://web.cs.ucla.edu/~malzantot/

http://www.planetb.ca/syntax-highlight-word