

Optimal Policy and Value Iteration:

a)

①: ① We choose to stay.

$$\begin{aligned}\sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t) &= \gamma^0 (-1) + \gamma^1 (-1) + \gamma^2 (-1) + \dots \\ &= -1 [ 1 + \gamma + \gamma^2 + \gamma^3 + \dots ] \\ &= -1 \left[ \frac{1}{1-\gamma} \right] \quad \left[ \because \text{sum of the geometric series here} \right] \\ &= \frac{-1}{1-\gamma}\end{aligned}$$

b)

② optimal actions would be to go from  $s_1$  to  $s_2$  and then come back to state  $s_1$ .

$$(s_1, s_2) \Rightarrow (go, go)$$

$$\begin{aligned}\text{Discounted Rewards} &= -2\gamma^0 + 3\gamma^1 \\ &= -2\gamma^0 + 3\gamma^1\end{aligned}$$

c)

$$V^0 = [0, 0] \quad V^1(s_1) = 0, V^1(s_2) = 0, T = 1$$

$$P(s_1 | s_1, a = go) = 1$$

$$P(s_1 | s_2, a = stay) = 0$$

$$P(s_2 | s_1, a = go) = 1$$

$$P(s_2 | s_1, a = stay) = 0$$

$$P(s_2 | s_2, a = stay) = 1$$

$$P(s_1 | s_1, a = stay) = 1$$

$$P(s_1 | s_1, a = go) = 0$$

$$P(s_2 | s_2, a = go) = 0$$

$$V^1 = V^1(s_1) = \max_a \left[ P(s_1 | s_1, a = go) [r(s_1, a = go) + \gamma V^0(s_1)], \right. \\ \left. P(s_1 | s_1, a = stay) [r(s_1, a = stay) + \gamma V^0(s_1)] \right]$$

$$= \max [ 1 [(-2) + 1 \times (0)], 1 [(-1) + 1 \times (0)] ]$$

$$= \max [ -2, -1 ]$$

$$= -1$$

$$V^1(s_2) = \max [ r(s_2, stay) + \gamma V^0(s_2), r(s_2, go) ]$$

$$= [ (-1 + 0), 3 ]$$

$$= 3$$

$$V^1 = [-1, 3]$$

Similarly, we can calculate for other

$$V^2 = V^2(s_1) = \max \left[ P(s_1 | s_1, a = go) [r(s_1, a = go) + \gamma V^1(s_1)], \right. \\ \left. P(s_1 | s_1, a = stay) [r(s_1, a = stay) + \gamma V^1(s_1)] \right]$$

&lt;

$$= \max [ 1 [ (-2) + 1 (3) ], \\ 1 [ -1 + (-1) ] ]$$

$$= \max [ 1, -2 ]$$

$$= 1$$

$$V^2(s_2) = \max (r(s_2, \text{stay}) + 1 V^1(s_2), r(s_2, \text{go}))]$$

$$\max ( -1 + 1 (3), 3 ]$$

$$\max ( 2, 3 )$$

$$3$$

$$V^2 = [ 1, 3 ]$$

$$V^3 = V^3(s_1) = \max [ P(s_1 | s_2, a = \text{go}) [r(s_2, a = \text{go}) + 1 V^2(s_2)], \\ P(s_1, s_2, a = \text{stay}) [r(s_1, a = \text{stay}) + V$$

$$= \max [ 1 [ (-2) + 1 (3) ],$$

$$1 [ (-1) + 1 (-2) ] ]$$

$$= \max ( 1, -3 )$$

$$V^3(s_1) = 1$$

$$V^3(s_2) = \max ( r(s_2, \text{stay}) + 1 V^2(s_2), \\ r(s_2, \text{go})$$

$$= \max ( (-1) + 3 ),$$

$$3 )$$

$$= 3$$

$$V^3 = [ 1, 3 ]$$

&lt;

$$\begin{aligned}
 V^2(s_2) &= \max [r(s_2, \text{stay}) + 1 V^1(s_2), r(s_2, \text{go})] \\
 &= \max (-1 + 1(3), 3) \\
 &= \max (2, 3) \\
 &= 3
 \end{aligned}$$

$$V^2 = [1, 3]$$

$$\begin{aligned}
 V^3(s_1) &= \max [p(s_1, s_2, a=\text{go}) [r(s_2, a=\text{go}) + 1 V^2(s_2)], \\
 &\quad p(s_1, s_2, a=\text{stay}) [r(s_1, a=\text{stay}) + V^2(s_1)]] \\
 &= \max [1 [(-2) + 1(3)], \\
 &\quad 1 [(-1) + 1(-2)]] \\
 &= \max (1, -3)
 \end{aligned}$$

$$V^3(s_1) = 1$$

$$\begin{aligned}
 V^3(s_2) &= \max (r(s_2, \text{stay}) + 1 V^2(s_2), \\
 &\quad r(s_2, \text{go})) \\
 &= \max ((-1) + 3, 3) \\
 &= 3
 \end{aligned}$$

$$V^3 = [1, 3]$$

So, Our answer:

$$V^1 = [-1, 3]$$

$$V^2 = [1, 3]$$

$$V^3 = [1, 3]$$

## 2: Value Iteration Convergence

a)

### ②. Value Iteration Convergence

$$\textcircled{a}. \quad V^1 = [-1, 3]$$

$$V^2 = [1, 3]$$

$$V^3 = [1, 3]$$

Since our value iteration does not change  
in  $V^2, V^3$ , we can say that

$$V^* = [1, 3]$$

$$\begin{aligned} \max_{s \in S} |V^1(s) - V^*(s)| &= \max_{s \in S} |V^1 - V^*| \\ &= \max_{s \in S} |[-1, 3] - [1, 3]| \\ &= \max_{s \in S} |(2), (0)| \\ &= 2 \end{aligned}$$

$$\begin{aligned} \max |V^2 - V^*| &= \max |[1, 3] - [1, 3]| \\ &= 0 \end{aligned}$$

$$\begin{aligned} \max |V^3 - V^*| &= \max |[1, 3] - [1, 3]| \\ &= 0 \end{aligned}$$

We can see from above that the error  
is decreasing monotonically.

b)

⑥: To prove:  $\|T(V) - T(V')\|_\infty \leq \gamma \|V - V'\|_\infty$

for any  $V, V' \in \mathbb{R}^{1^S}$

$$\begin{aligned} \|T(V) - T(V')\|_\infty &= \max_{s \in S} |T(V)(s) - T(V')(s)| \\ &= \max_{s \in S} \left| \max_{a \in A} \sum_{s' \in S} P(s, a, s') (R(s, a) + \gamma V(s')) - \max_{a \in A} \sum_{s' \in S} P(s, a, s') (R(s, a) + \gamma V'(s')) \right| \end{aligned}$$

To continue, we will look at some simple inequality expressions and then evaluate the above eq<sup>n</sup>.

Let's suppose  $f: X \rightarrow \mathbb{R}$ , and  $g: S \rightarrow \mathbb{R}$   
then

$$\forall x, f(x) - g(x) \leq |f(x) - g(x)|$$

$$\forall x, f(x) \leq |f(x) - g(x)| + g(x)$$

$$\max_{x \in X} f(x) \leq \max_{x \in X} |f(x) - g(x)| + \max_{x \in X} g(x)$$

$$\max_{x \in X} f(x) - \max_{x \in X} g(x) \leq \max_{x \in X} |f(x) - g(x)|$$

Now, if  $\max_{x \in X} f(x) - \max_{x \in X} g(x) \geq 0$ , then

$$\left| \max_{x \in X} f(x) - \max_{x \in X} g(x) \right| \leq \max_{x \in X} |f(x) - g(x)| \quad \text{--- (2)}$$

Using the above eq<sup>n</sup> we can rewrite eq<sup>n</sup> ⑥ as

$$\begin{aligned} \|T(V) - T(V')\| &\leq \max_{s \in S} \max_{a \in A} \left| \sum_{s' \in S} P(s, a, s') (R(s, a) + \gamma V(s')) - \sum_{s' \in S} P(s, a, s') (R(s, a) + \gamma V'(s')) \right| \\ &= \gamma \max_{s \in S} \max_{a \in A} \left| \sum_{s' \in S} P(s, a, s') (V(s') - V'(s')) \right| \\ &= \gamma \max_{s \in S} \max_{a \in A} \sum_{s' \in S} P(s, a, s') |V(s') - V'(s')| \end{aligned}$$

Using eq<sup>n</sup> ③.

$$\begin{aligned} &\leq \gamma \max_{s \in S} \max_{a \in A} \max_{s' \in S} |V(s') - V'(s')| \\ &= \gamma \max_{s' \in S} |V(s') - V'(s')| \\ &= \gamma \|V(s') - V'(s')\| \end{aligned}$$

Thus, we can say that Bellman operator is a contraction mapping, and value iteration function converges to a unique fixed point.

c)

$$\textcircled{C}: \|v^{n+1} - v^*\|_\infty = \|v^{n+1} - T(v^{n+1}) + T(v^{n+1}) - v^*\|$$

Now we can use triangle inequality here:

$$\|a+b\| \leq \|a\| + \|b\|$$

$$\begin{aligned} \|v^{n+1} - v^*\| &\leq \|v^{n+1} - T(v^{n+1})\| + \|T(v^{n+1}) - v^*\| \\ &\leq \|T(v^n) - T(v^{n+1})\| + \|T(v^{n+1}) - T(v^*)\| \\ &\quad [\because v^{n+1} = T(v^n) \text{ \& } T(v^*) = v^*] \end{aligned}$$

$$\begin{aligned} &\leq \gamma \|v^n - v^{n+1}\| + \gamma \|v^{n+1} - v^*\| \\ &\quad [\because \|T(v) - T(v')\| \leq \gamma \|v - v'\| \text{ as proved} \\ &\quad \text{in part B of the assignment}] \end{aligned}$$

$$(1-\gamma) \|v^{n+1} - v^*\| \leq \gamma \|v^n - v^{n+1}\|$$

$$\|v^{n+1} - v^*\| \leq \frac{\gamma}{1-\gamma} \|v^n - v^{n+1}\|$$

Now as we know that  $\{v^n\}$  sequence has the property of being a Cauchy sequence, which essentially means

$$\|v^n - v^{n+1}\| < \epsilon$$

where  $\epsilon \rightarrow 0$  [A small positive No.]

$$\text{Therefore } \|v^{n+1} - v^*\| \leq \frac{\gamma}{1-\gamma} \epsilon$$

d)- BONUS

<

⑧: Given  $\|T(x_1) - T(x_2)\|_\infty \leq \alpha \|x_1 - x_2\|_\infty$ ,  
 $T$  has a fixed point,  $\exists x^*: T(x^*) = x^*$   
 fixed point is unique.

Consider the sequence  $\{v_i\}$  where  $v_i = T v_{i-1}$   
 beginning with  $v_0$ . So the sequence would  
 be like  $v, T v, T^2 v, \dots$

So, from given eq<sup>n</sup>

$$\|T v - T^2 v\| \leq \alpha \|v - T v\|$$

$$\|T^2 v - T^3 v\| \leq \alpha \|T v - T^2 v\|$$

$$\leq \alpha^2 \|v - T v\|$$

⋮

$$\|T^R v - T^{R+1} v\| \leq \alpha^R \|v - T v\|$$

Now since  $\alpha < 1$ , and the sequence  
 $v, T v, T^2 v, \dots$  is a Cauchy sequence and  
 as it is over the Euclidean space, the  
 sequence should converge.

As the sequence converges, there exist some  
 final value  $v$  such that  $\bar{v} = T \bar{v}$ .

Why is it unique?

Suppose  $\hat{v} \neq \bar{v}$  is a fixed point of  $T$ .

We then have  $\hat{v} = T \hat{v}$  and

$$\|\bar{v} - \hat{v}\| = \|T \bar{v} - T \hat{v}\|$$

$$\leq \alpha \|\bar{v} - \hat{v}\|$$

and since  $\alpha < 1$ , we must have  $\bar{v} = \hat{v}$

and therefore which contradicts our initial  
 statement of  $\hat{v} \neq \bar{v}$ , hence  $\hat{v}$  is unique.



#### 4: Policy Gradient Variance Reduction

a)

$$(4) \quad \nabla_{\theta} J(\theta) = \nabla_{\theta} E_{\tau \sim \pi_{\theta}} [R(\tau)]$$

$$\approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log \pi_{\theta}(\tau_i)$$

$$\text{if } R(\tau) = R(s) - b$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} E_{\tau \sim \pi_{\theta}} [R(\tau)] \\ &= E_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left( \sum_{t'=t}^{T-1} r_{t'} \right) - \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t) \right] \end{aligned}$$

We need to show that for any single time  $t$ , the gradient of  $\log \pi_{\theta}(a_t | s_t)$  multiplied by  $b(s_t)$  is zero.

$$\begin{aligned} E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)] &= E_{s_{0:t}, a_{0:t-1}} [E_{s_{t+1:T}, a_{t:T-1}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)]] \\ &= E_{s_{0:t}, a_{0:t-1}} [E_{s_{t+1:T}, a_{t:T-1}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)]] \\ &= E_{s_{0:t}, a_{0:t-1}} [E_{s_{t+1:T}, a_{t:T-1}} [b(s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]] \end{aligned}$$

Note that I have written trajectory  $s_0, a_0, \dots, a_{T-1}, s_T$  as  $s_{0:T}, a_{0:T-1}$ .

Now,

$$E_{a_t} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] = 0, \text{ because we can write.}$$

$$\begin{aligned} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) &= \int \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \pi_{\theta}(a_t | s_t) da_t \\ &= \nabla_{\theta} \int \pi_{\theta}(a_t | s_t) da_t \\ &= \nabla_{\theta} \cdot 1 \\ &= 0 \end{aligned}$$

$$\text{Hence } E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)] = 0$$

This shows that  $b$  does not cause bias.

b)

$$b) \quad \text{Var}(x) = E(x^2) - (E(x))^2$$

$$x = R(C) \nabla_{\theta} \log \pi_{\theta}(C)$$

with baseline

$$x = (R(C) - b) \nabla_{\theta} \log \pi_{\theta}(C)$$

$$V(x) = E((R(C) - b)^2 \nabla_{\theta} \log \pi_{\theta}(C)^2) - (E((R(C) - b) \nabla_{\theta} \log \pi_{\theta}(C)))^2$$

We already showed in part a), that introducing a baseline does not cause bias, hence  $(E(x))^2 \approx 0$ , so we are left with the  $E(x^2)$  term.

$$V(x) = \sum_{t=0}^T E_t \left[ \left( \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (R(C) - b) \right)^2 \right]$$

We know that if 2 variables are independent,

then  $E(xy) = E(x)E(y)$ , similarly in above due to independence we can write:

$$V(x) = \sum_{t=0}^T E_t \left[ \left( \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right)^2 \right] E_t \left[ (R(C) - b)^2 \right]$$

Here if we are able to optimize our choice of  $b(s_t)$

then  $E_t[(R(C) - b(s_t))^2]$  is a least square problem.

and the minimum value of this expectation can

occur when  $b(s_t) \approx R(C)$  to approximate the expected return starting at time  $t$ .

# dp

November 10, 2019

## 1 Dynamic Programming (20 points + 10 bonus points)

In this assignment, we will implement a few dynamic programming algorithms, namely, policy iteration and value iteration and run them on a simple MDP - the Frozen Lake environment.

The sub-routines for these algorithms are present in `vi_and_pi.py` and must be filled out to test your implementation.

The deliverables are located at the end of this notebook and show the point distribution for each part.

**Value iteration is worth 20 points of regular credit and policy iteration is worth 10 points of bonus credit for both sections of this course CS 7643 and CS 4803.**

```
[67]: %load_ext autoreload
      %autoreload 2

      import numpy as np
      import gym
      import time

      from IPython.display import clear_output

      from lake_envs import *
      from vi_and_pi import *

      np.set_printoptions(precision=3)

      env_d = gym.make("Deterministic-4x4-FrozenLake-v0")
      env_s = gym.make("Stochastic-4x4-FrozenLake-v0")
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

### 1.0.1 Render Mode

The variable `RENDER_ENV` is set `True` by default to allow you to see a rendering of the state of the environment at every time step. However, when you complete this assignment, you must set this to `False` and re-run all blocks of code. This is to prevent excessive amounts of rendered environments from being included in the final PDF.

**IMPORTANT: SET RENDER\_ENV TO FALSE BEFORE SUBMISSION!**

```
[126]: RENDER_ENV = False
```

## 1.1 Part 1: Value Iteration

For the first part, you will implement the familiar value iteration update from class.

In `vi_and_pi.pi` and complete the `value_iteration` function.

```
[127]: #####
# Use this space for debugging                                #
# Make sure to delete this code before submission #
#####
pass
#####
```

Run the cell below to train value iteration and render a single episode of following the policy obtained at the end of value iteration.

You should expect to get an Episode reward of 1.0.

```
[128]: print("\n" + "-"*25 + "\nBeginning Value Iteration\n" + "-"*25)

V_vi, p_vi = value_iteration(env_d.P, env_d.nS, env_d.nA, gamma=0.9, tol=1e-3)
render_single(env_d, p_vi, 100, show_rendering=RENDER_ENV)
```

```
-----
Beginning Value Iteration
-----
```

```
Episode reward: 1.000000
```

## 1.2 [BONUS] Part 2: Policy Iteration

This is a bonus question in which you will implement policy iteration. If you do not wish to attempt this bonus question, skip to the next part.

In class, we studied the value iteration update:

$$V_{t+1}(s) \leftarrow \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V_t(s')]$$

This is used to compute the value function  $V^*$  corresponding to the optimal policy  $\pi^*$ . We can alternatively compute the value function  $V^\pi$  corresponding to an arbitrary policy  $\pi$ , with a similar update loop:

$$V_{t+1}^\pi(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V_t^\pi(s')]$$

On convergence, this will give us  $V^\pi$ , which is the first step of a policy iteration update.

The second step involves policy refinement, which will update the policy to take actions greedily with respect to  $V^\pi$ :

$$\pi_{new} \leftarrow \arg \max_a \left[ r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^\pi(s') \right]$$

A single update of policy iteration involves the two above steps: (1) policy evaluation (which itself is an inner loop which will converge to  $V^\pi$  and (2) policy refinement. In the first part of assignment, you will implement the functions for policy evaluation, policy improvement (refinement) and policy iteration.

In `vi_and_pi.pi` and complete the `policy_evaluation`, `policy_improvement` and `policy_iteration` functions. Run the blocks below to test your algorithm.

```
[129]: #####
# Use this space for debugging                                #
# Make sure to delete this code before submission #
#####
pass
#####
```

```
[130]: print("\n" + "-"*25 + "\nBeginning Policy Iteration\n" + "-"*25)

V_pi, p_pi = policy_iteration(env_d.P, env_d.nS, env_d.nA, gamma=0.9, tol=1e-3)
render_single(env_d, p_pi, 100, show_rendering=RENDER_ENV)
```

```
-----
Beginning Policy Iteration
-----
Episode reward: 1.000000
```

### 1.3 Part 3: VI on Stochastic Frozen Lake

Now we will apply our implementation on an MDP where transitions to next states are stochastic. Modify your implementation of value iteration as needed so that policy iteration and value iteration work for stochastic transitions.

```
[131]: #####
# Use this space for debugging                                #
# Make sure to delete this code before submission #
#####
pass
#####
```

```
[134]: print("\n" + "-"*25 + "\nBeginning Value Iteration\n" + "-"*25)

V_vi, p_vi = value_iteration(env_s.P, env_s.nS, env_s.nA, gamma=0.9, tol=1e-3)
```

```
render_single(env_s, p_vi, 100, show_rendering=RENDER_ENV)
```

```
-----  
Beginning Value Iteration  
-----
```

```
Episode reward: 1.000000
```

## 1.4 [BONUS] Part 4: PI on Stochastic Frozen Lake

This is a bonus question to run policy iteration on stochastic frozen lake.

Now we will apply our implementation on an MDP where transitions to next states are stochastic. Modify your implementation of value iteration as needed so that policy iteration and value iteration work for stochastic transitions.

```
[135]: #####  
      # Use this space for debugging                                #  
      # Make sure to delete this code before submission #  
      #####  
      pass  
      #####
```

```
[136]: print("\n" + "-"*25 + "\nBeginning Policy Iteration\n" + "-"*25)  
  
V_pi, p_pi = policy_iteration(env_s.P, env_s.nS, env_s.nA, gamma=0.9, tol=1e-3)  
render_single(env_s, p_pi, 100, show_rendering=RENDER_ENV)
```

```
-----  
Beginning Policy Iteration  
-----
```

```
Episode reward: 1.000000
```

## 1.5 Evaluate All Policies

Now, we will first test the value iteration implementation on two kinds of environments - the deterministic FrozenLake and the stochastic FrozenLake. We will also run the same for policy iteration

### 1.5.1 Deliverable 1 (10 points)

Run value iteration on deterministic FrozenLake. You should get a reward of 1.0 for full credit.

```
[137]: print("\nValue Iteration on Deterministic FrozenLake:")  
V_vi, p_vi = value_iteration(env_d.P, env_d.nS, env_d.nA, gamma=0.9, tol=1e-3)  
evaluate(env_d, p_vi, max_steps=100, max_episodes=2)
```

```
Value Iteration on Deterministic FrozenLake:
```

```
> Average reward over 2 episodes: 1.0
> Percentage of episodes goal reached: 100%
```

### 1.5.2 Deliverable 2 (10 points)

Run value iteration on stochastic FrozenLake. Note that this time, running the same policy over multiple episodes will result in different outcomes (final reward) due to stochastic transitions in the environment, and even the optimal policy may not succeed in reaching the goal state 100% of the time.

You should get a reward of 0.7 or higher over 1000 episodes for full credit.

```
[138]: print("\nValue Iteration on Stochastic FrozenLake:")
V_vi, p_vi = value_iteration(env_s.P, env_s.nS, env_s.nA, gamma=0.9, tol=1e-3)
evaluate(env_s, p_vi, max_steps=100, max_episodes=1000)
```

```
Value Iteration on Stochastic FrozenLake:
> Average reward over 1000 episodes: 0.737
> Percentage of episodes goal reached: 94%
```

### 1.5.3 Deliverable 3 (5 bonus points)

Run policy iteration on deterministic FrozenLake. You should get a reward of 1.0 for full credit.

```
[139]: print("Policy Iteration on Deterministic FrozenLake:")
V_pi, p_pi = policy_iteration(env_d.P, env_d.nS, env_d.nA, gamma=0.9, tol=1e-3)
evaluate(env_d, p_pi, max_steps=100, max_episodes=2)
```

```
Policy Iteration on Deterministic FrozenLake:
> Average reward over 2 episodes: 1.0
> Percentage of episodes goal reached: 100%
```

### 1.5.4 Deliverable 4 (5 bonus points)

Run policy iteration on stochastic FrozenLake.

You should get a reward of 0.7 or higher over 1000 episodes for full credit.

```
[140]: print("Policy Iteration on Stochastic FrozenLake:")
V_pi, p_pi = policy_iteration(env_s.P, env_s.nS, env_s.nA, gamma=0.9, tol=1e-3)
evaluate(env_s, p_pi, max_steps=100, max_episodes=1000)
```

```
Policy Iteration on Stochastic FrozenLake:
> Average reward over 1000 episodes: 0.738
> Percentage of episodes goal reached: 93%
```

## 1.6 Submission Reminder

**PLEASE RE-RUN THE NOTEBOOK WITH `RENDER_ENV` SET TO `FALSE` BEFORE SUBMISSION!**

## Q-Learning & DQNs (30 points + 5 bonus points)

In this section, we will implement a few key parts of the Q-Learning algorithm for two cases - (1) A Q-network which is a single linear layer (referred to in RL literature as "Q-learning with linear function approximation") and (2) A deep (convolutional) Q-network, for some Atari game environments where the states are images.

Optional Readings:

- **Playing Atari with Deep Reinforcement Learning**, Mnih et. al., <https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf> (<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>)
- **The PyTorch DQN Tutorial** [https://pytorch.org/tutorials/intermediate/reinforcement\\_q\\_learning.html](https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html) ([https://pytorch.org/tutorials/intermediate/reinforcement\\_q\\_learning.html](https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html))

Note: The bonus credit for this question applies to both sections CS 7643 and CS 4803

In [5]:

```
!ls
```

```
core  sample_data
```

In [8]:

```
%load_ext autoreload
%autoreload 2

import numpy as np
import gym

import torch
import torch.nn as nn
import torch.optim as optim

from core.dqn_train import DQNTrain
from utils.test_env import EnvTest
from utils.schedule import LinearExploration, LinearSchedule
from utils.preprocess import greyscale
from utils.wrappers import PreproWrapper, MaxAndSkipEnv

from linear_qnet import LinearQNet
from cnn_qnet import ConvQNet

if torch.cuda.is_available():
    device = torch.device('cuda', 0)
else:
    device = torch.device('cpu')
```

The autoreload extension is already loaded. To reload it, use:  
%reload\_ext autoreload



## Part 1: Setup Q-Learning with Linear Function Approximation

Training Q-networks using (Deep) Q-learning involves a lot of moving parts. However, for this assignment, the scaffolding for the first 3 points listed below is provided in full and you must only complete point 4. You may skip to point 4 if you only care about the implementation required for this assignment.

1. **Environments:** We will use the standardized OpenAI Gym framework for environment API calls (read through <http://gym.openai.com/docs/> (<http://gym.openai.com/docs/>) if you want to know more details about this interface). Specifically, we will use a custom Test environment defined in `utils/test_env.py` for initial sanity checks and then Gym-Atari environments later on.
1. **Exploration:** In order to train any RL model, we require experience or "data" gathered from interacting with the environment by taking actions. What policy should we use to collect this experience? Given a Q-network, one may be tempted to define a greedy policy which always picks the highest valued action at every state. However, this strategy will in most cases not work since we may get stuck in a local minima and never explore new states in the environment which may lead to a better reward. Hence, for the purpose of gathering experience (or "data") from the environment, it is useful to follow a policy that deviates from the greedy policy slightly in order to explore new states. A common strategy used in RL is to follow an  $\epsilon$ -greedy policy which with probability  $0 < \epsilon < 1$  picks a random action instead of the action provided by the greedy policy.
1. **Replay Buffers:** Data gathered from a single trajectory of states and actions in the environment provides us with a batch of highly correlated (non IID) data, which leads to high variance in gradient updates and convergence. In order to ameliorate this, replay buffers are used to gather a set of transitions i.e. (state, action, reward, next state) tuples, by executing multiple trajectories in the environment. Now, for updating the Q-Network, we will first wait to fill up our replay buffer with a sufficiently large number of transitions over multiple different trajectories, and then randomly sample a batch of transitions to compute loss and update the models.
1. **Q-Learning network, loss and update:** Finally, we come to the part of Q-learning that we will implement for this assignment -- the Q-network, loss function and update. In particular, we will implement a variant of Q-Learning called "Double Q-Learning", where we will maintain two Q networks -- the first Q network is used to pick actions and the second "target" Q network is used to compute Q-values for the picked actions. Here is some reference material on the same - [Blog 1 \(https://towardsdatascience.com/double-q-learning-the-easy-way-a924c4085ec3\)](https://towardsdatascience.com/double-q-learning-the-easy-way-a924c4085ec3), [Blog 2 \(https://medium.com/@ameetsd97/deep-double-q-learning-why-you-should-use-it-bedf660d5295\)](https://medium.com/@ameetsd97/deep-double-q-learning-why-you-should-use-it-bedf660d5295), but we will not need to get into the details of Double Q-learning for this assignment. Now, let's walk through the steps required to implement this below.
  - **Linear Q-Network:** In `linear_qnet.py`, define the initialization and forward pass of a Q-network with a single linear layer which takes the state as input and outputs the Q-values for all actions.
  - **Setting up Q-Learning:** In `core/dqn_train.py`, complete the functions `process_state`, `forward_loss` and `update_step` and `update_target_params`. The loss function for our Q-Networks is defined for a single transition tuple of (state, action, reward, next state) as follows.  $Q(s_t, a_t)$  refers to the state-action values computed by our first Q-network at the current state and for the current actions,  $Q_{target}(s_{t+1}, a_{t+1})$  refers to the state-action values for the next state and all possible future actions computed by the target Q-Network

$$\begin{aligned}
 Q_{sample}(s_t) &= r_t \text{ if done} \\
 &= r_t + \gamma \max_{a_{t+1}} Q_{target}(s_{t+1}, a_{t+1}) \text{ otherwise} \\
 \text{Loss} &= (Q(s_t, a_t) - Q(s_{t+1}, a_{t+1}))^2
 \end{aligned}$$

## Deliverable 1 (15 points)

Run the following block of code to train a Linear Q-Network. You should get an average reward of ~4.0, full credit will be given if average reward at the final evaluation is above 3.5

```
In [10]: from configs.pl_linear import config as config_lin

env = EnvTest((5, 5, 1))

# exploration strategy
exp_schedule = LinearExploration(env, config_lin.eps_begin,
                                config_lin.eps_end, config_lin.eps_nsteps)

# learning rate schedule
lr_schedule = LinearSchedule(config_lin.lr_begin, config_lin.lr_end,
                              config_lin.lr_nsteps)

# train model
model = DQNTrain(LinearQNet, env, config_lin, device)
model.run(exp_schedule, lr_schedule)
```

Evaluating...

Average reward: 2.30 +/- 0.00

1001/10000 [==>.....] - ETA: 7s - Loss: 4.0490 -  
Avg\_R: 0.9300 - Max\_R: 3.1000 - eps: 0.8020 - Grads: 14.9250 - Max\_Q:  
0.8754 - lr: 0.0042

Evaluating...

Average reward: 3.90 +/- 0.00

2001/10000 [=====>.....] - ETA: 7s - Loss: 11.4660  
- Avg\_R: 1.1300 - Max\_R: 4.1000 - eps: 0.6040 - Grads: 21.6438 - Max\_Q:  
1.9417 - lr: 0.0034

Evaluating...

Average reward: 3.90 +/- 0.00

3001/10000 [=====>.....] - ETA: 6s - Loss: 7.5043 -  
Avg\_R: 2.3600 - Max\_R: 4.1000 - eps: 0.4060 - Grads: 31.8965 - Max\_Q:  
2.3320 - lr: 0.0026

Evaluating...

Average reward: 3.80 +/- 0.00

4001/10000 [=====>.....] - ETA: 5s - Loss: 4.4914 -  
Avg\_R: 3.2600 - Max\_R: 4.1000 - eps: 0.2080 - Grads: 12.3865 - Max\_Q:  
2.4202 - lr: 0.0018

Evaluating...

Average reward: 3.80 +/- 0.00

5001/10000 [=====>.....] - ETA: 4s - Loss: 2.6182 -  
Avg\_R: 3.8550 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 32.6568 - Max\_Q:  
2.7201 - lr: 0.0010

Evaluating...

Average reward: 4.10 +/- 0.00

6001/10000 [=====>.....] - ETA: 3s - Loss: 0.1513 -  
Avg\_R: 4.1000 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 2.4042 - Max\_Q: 2.  
4846 - lr: 0.0010

Evaluating...

Average reward: 4.10 +/- 0.00

7001/10000 [=====>.....] - ETA: 2s - Loss: 0.7107 -  
Avg\_R: 4.0950 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 7.2632 - Max\_Q: 2.  
5210 - lr: 0.0010

Evaluating...

Average reward: 3.80 +/- 0.00

```
8001/10000 [=====>.....] - ETA: 1s - Loss: 0.0104 -
Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 1.9160 - Max_Q: 2.
8098 - lr: 0.0010
```

```
Evaluating...
Average reward: 4.10 +/- 0.00
```

```
9001/10000 [=====>...] - ETA: 0s - Loss: 1.7742 -
Avg_R: 3.9050 - Max_R: 4.0000 - eps: 0.0100 - Grads: 17.5330 - Max_Q:
2.8558 - lr: 0.0010
```

```
Evaluating...
Average reward: 3.90 +/- 0.00
```

```
10001/10000 [=====] - 9s - Loss: 0.7037 - Avg_
R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 11.9553 - Max_Q: 2.878
4 - lr: 0.0010
```

```
- Training done.
Evaluating...
```

```
Average reward: 4.10 +/- 0.00
```

You should get a final average reward of over 4.0 on the test environment.

## Part 2: Q-Learning with Deep Q-Networks

In `cnn_qnet.py`, implement the initialization and forward pass of a convolutional Q-network with architecture as described in this DeepMind paper:

"Playing Atari with Deep Reinforcement Learning", Mnih et. al.

(<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf> (<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>))

### Deliverable 2 (10 points)

Run the following block of code to train our Deep Q-Network. You should get an average reward of ~4.0, full credit will be given if average reward at the final evaluation is above 3.5

```
In [13]: from configs.p2_cnn import config as config_cnn

env = EnvTest((80, 80, 1))

# exploration strategy
exp_schedule = LinearExploration(env, config_cnn.eps_begin,
                                  config_cnn.eps_end, config_cnn.eps_nsteps)

# learning rate schedule
lr_schedule = LinearSchedule(config_cnn.lr_begin, config_cnn.lr_end,
                              config_cnn.lr_nsteps)

# train model
model = DQNTrain(ConvQNet, env, config_cnn, device)
model.run(exp_schedule, lr_schedule)
```

Evaluating...

Average reward: 0.00 +/- 0.00

Populating the memory 150/200...

Evaluating...

Average reward: 0.00 +/- 0.00

301/1000 [=====>.....] - ETA: 1s - Loss: 0.7483 - Avg\_R: -0.0750 - Max\_R: 2.3000 - eps: 0.4060 - Grads: 13.9815 - Max\_Q: 0.0951 - lr: 0.0002

Evaluating...

Average reward: -1.00 +/- 0.00

401/1000 [=====>.....] - ETA: 1s - Loss: 4.3078 - Avg\_R: -0.2350 - Max\_R: 3.8000 - eps: 0.2080 - Grads: 86.2777 - Max\_Q: 0.1893 - lr: 0.0001

Evaluating...

Average reward: 0.50 +/- 0.00

501/1000 [=====>.....] - ETA: 1s - Loss: 3.9153 - Avg\_R: 0.4850 - Max\_R: 2.3000 - eps: 0.0100 - Grads: 92.5900 - Max\_Q: 0.1937 - lr: 0.0001

Evaluating...

Average reward: 0.50 +/- 0.00

601/1000 [=====>.....] - ETA: 1s - Loss: 5.5912 - Avg\_R: 2.3050 - Max\_R: 4.0000 - eps: 0.0100 - Grads: 102.4696 - Max\_Q: 0.2438 - lr: 0.0001

Evaluating...

Average reward: 4.00 +/- 0.00

701/1000 [=====>.....] - ETA: 0s - Loss: 3.4055 - Avg\_R: 4.0300 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 53.8625 - Max\_Q: 0.3494 - lr: 0.0001

Evaluating...

Average reward: 4.10 +/- 0.00

801/1000 [=====>.....] - ETA: 0s - Loss: 1.8910 - Avg\_R: 4.0550 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 35.5276 - Max\_Q: 0.4453 - lr: 0.0001

Evaluating...

Average reward: 4.10 +/- 0.00

901/1000 [=====>...] - ETA: 0s - Loss: 1.4227 - Avg\_R: 3.6000 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 51.5050 - Max\_Q: 0.5274 - lr: 0.0001

Evaluating...

Average reward: 4.10 +/- 0.00

1001/1000 [=====>] - 3s - Loss: 0.6211 - Avg\_R: 4.0700 - Max\_R: 4.1000 - eps: 0.0100 - Grads: 76.9530 - Max\_Q: 0.5990 - lr: 0.0001

```
- Training done.  
Evaluating...  
Average reward: 4.10 +/- 0.00
```

You should get a final average reward of over 4.0 on the test environment, similar to the previous case.

## Part 3: Playing Atari Games from Pixels - using Linear Function Approximation

Now that we have setup our Q-Learning algorithm and tested it on a simple test environment, we will shift to a harder environment - an Atari 2600 game from OpenAI Gym: Pong-v0 (<https://gym.openai.com/envs/Pong-v0/>), where we will use RGB images of the game screen as our observations for state.

No additional implementation is required for this part, just run the block of code below (will take around 1 hour to train). We don't expect a simple linear Q-network to do well on such a hard environment - full credit will be given simply for running the training to completion irrespective of the final average reward obtained.

You may edit `configs/p3_train_atari_linear.py` if you wish to play around with hyperparameters for improving performance of the linear Q-network on Pong-v0, or try another Atari environment by changing the `env_name` hyperparameter. The list of all Gym Atari environments are available here: <https://gym.openai.com/envs/#atari>

### Deliverable 3 (5 points)

Run the following block of code to train a linear Q-network on Atari Pong-v0. We don't expect the linear Q-Network to learn anything meaningful so full credit will be given for simply running this training to completion (without errors), irrespective of the final average reward.



```
In [12]: from configs.p3_train_atari_linear import config as config_lina

# make env
env = gym.make(config_lina.env_name)
env = MaxAndSkipEnv(env, skip=config_lina.skip_frame)
env = PreproWrapper(env, prepro=greyscale, shape=(80, 80, 1),
                    overwrite_render=config_lina.overwrite_render)

# exploration strategy
exp_schedule = LinearExploration(env, config_lina.eps_begin,
                                config_lina.eps_end, config_lina.eps_nsteps)

# learning rate schedule
lr_schedule = LinearSchedule(config_lina.lr_begin, config_lina.lr_end,
                             config_lina.lr_nsteps)

# train model
model = DQNTrain(LinearQNet, env, config_lina, device)
print("Linear Q-Net Architecture:\n", model.q_net)
model.run(exp_schedule, lr_schedule)
```

Evaluating...

Linear Q-Net Architecture:

```
LinearQNet(
  (fully_connected): Linear(in_features=25600, out_features=6, bias=True)
)
```

Average reward: -20.80 +/- 0.06

250301/500000 [=====>.....] - ETA: 1208s - Loss: 4.9595 - Avg\_R: -20.4400 - Max\_R: -17.0000 - eps: 0.7747 - Grads: 496.0526 - Max\_Q: 5.5334 - lr: 0.0001

Evaluating...

Average reward: -21.00 +/- 0.00

500001/500000 [=====] - 2554s - Loss: 61.1647 - Avg\_R: -20.5200 - Max\_R: -19.0000 - eps: 0.5500 - Grads: 894.1898 - Max\_Q: 7.2395 - lr: 0.0001

- Training done.

Evaluating...

Average reward: -20.92 +/- 0.05

## Part 4: [BONUS] Playing Atari Games from Pixels - using Deep Q-Networks

This part is extra credit and worth 5 bonus points. We will now train our deep Q-Network from Part 2 on Pong-v0.

Again, no additional implementation is required but you may wish to tweak your CNN architecture in `cnn_qnet.py` and hyperparameters in `configs/p4_train_atari_cnn.py` (however, evaluation will be considered at no farther than the default 5 million steps, so you are not allowed to train for longer). Please note that this training may take a very long time (we tested this on a single GPU and it took around 6 hours).

The bonus points for this question will be allotted based on the best evaluation average reward (EAR) before 5 million time steps:

1. EAR  $\geq$  0.0 : 4/4 points
2. EAR  $\geq$  -5.0 : 3/4 points
3. EAR  $\geq$  -10.0 : 3/4 points
4. EAR  $\geq$  -15.0 : 1/4 points

### Deliverable 4: (5 bonus points)

Run the following block of code to train your DQN:

```
In [0]: from configs.p4_train_atari_cnn import config as config_cnn

# make env
env = gym.make(config_cnn.env_name)
env = MaxAndSkipEnv(env, skip=config_cnn.skip_frame)
env = PreproWrapper(env, prepro=greyscale, shape=(80, 80, 1),
                    overwrite_render=config_cnn.overwrite_render)

# exploration strategy
exp_schedule = LinearExploration(env, config_cnn.eps_begin,
                                config_cnn.eps_end, config_cnn.eps_nsteps)

# learning rate schedule
lr_schedule = LinearSchedule(config_cnn.lr_begin, config_cnn.lr_end,
                             config_cnn.lr_nsteps)

# train model
model = DQNTrain(ConvQNet, env, config_cnn, device)
print("CNN Q-Net Architecture:\n", model.q_net)
model.run(exp_schedule, lr_schedule)
```