

CS7643: Deep Learning

Fall 2019

HW4 Solutions

James Hahn

November 12, 2019

1 Optimal Policy and Value Function

1. First, let's take the sum of discounted rewards as $\sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t)$

We know we start at S_1 , so $s_0 = S_1$, and the we always choose “stay”, so $a_0 = a_i = \text{“stay”}$. Since we always “stay” at the same state, s_i will always be S_1 . As such, this summation becomes:

$$\begin{aligned} & \sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t) \\ &= \sum_{t=0}^{\infty} \gamma^t r_t(S_0, \text{“stay”}) \\ &= \sum_{t=0}^{\infty} \gamma^t (-1) \\ &= - \sum_{t=0}^{\infty} \gamma^t \\ &= -\frac{1}{1-\gamma} \end{aligned}$$

So, if we assume this simulation runs for an infinite number of steps, then the sum of discounted rewards is the value $-\infty \cdot \gamma$, which becomes $-\infty$ since $\gamma > 0$.

2. First, let's observe how the value of γ changes the policy values. If we say, $\gamma = 0$, then this doesn't give any incentive to future rewards. If $\gamma = 1$, then this heavily incentivizes future rewards (i.e. reaching the goal state). The optimal policy depends on the value of γ . As such, the optimal policy when $\gamma \geq 0.232$, we want to choose the policy $(a_1, a_2) = (\text{"go"}, \text{"go"})$. Otherwise, when $\gamma < 0.232$, we want the policy $(a_1, a_2) = (\text{"stay"}, \text{"stay"})$. The sum of discounted rewards can be found below for both options:

Assume we start at S_1 with policy $(a_1, a_2) = (\text{"stay"}, \text{"stay"})$ when $\gamma < 0.232$. This results in a sum of discounted rewards of $\sum_{t=0}^{\infty} (-1)\gamma^t = -\frac{1}{1-\gamma}$. So, when $\gamma = 0.232$, we get the final reward as -1.303. Now, when $\gamma = 0$, we get the final reward as -1. This means this case's sum of discounted rewards is bounded between -1 and -1.303.

Assume we start at S_1 and use the policy $(a_1, a_2) = (\text{"go"}, \text{"go"})$ when $\gamma \geq 0.232$. This is optimal because a high γ value incentivizes future rewards all other rewards in the MDP are negative, except for the termination reward, which is a reward of +3. Assuming we're forced to take an action at each iteration of the simulation, the only way to achieve that positive reward is to first traverse to S_2 and then terminate the program by choosing the "go" action. This results in a sum of discounted rewards of $\sum_{t=0}^1 \gamma^t r_t(s_t, a_t) = \gamma^0 r_0(S_1, \text{"go"}) + \gamma^1 r_1(S_2, \text{"go"}) = -2 + \gamma(3) = 3\gamma - 2$. So, when $\gamma = 0.232$, we get the final reward as -1.303. Now, when $\gamma = 1$, we get the final reward as +1. This means this case's sum of discounted rewards is bounded between -1.303 and +1.

As shown above, the cutoff point between the two action policies is at $\gamma = 0.232$, and the optimal policies are provided for the over and under cases.

3. $V_0 = [0, 0]$

$$V_1 = [\max(r(s_1, \text{"stay"}) + \gamma V_0(s_1), r(s_1, \text{"go"}) + \gamma V_0(s_2)), \max(r(s_2, \text{"stay"}) + \gamma V_0(s_2), r(s_2, \text{"go"}))] = [\max(-1, -2), \max(-1, 3)] = [-1, 3]$$

$$V_2 = [\max(r(s_1, \text{"stay"}) + \gamma V_1(s_1), r(s_1, \text{"go"}) + \gamma V_1(s_2)), \max(r(s_2, \text{"stay"}) + \gamma V_1(s_2), r(s_2, \text{"go"}))] = [\max(-1 - \gamma, -2 + 3\gamma), \max(-1 + 3\gamma, 3)] = [\max(-2, 1), \max(2, 3)] = [1, 3]$$

$$V_3 = [\max(r(s_1, \text{"stay"}) + \gamma V_2(s_1), r(s_1, \text{"go"}) + \gamma V_2(s_2)), \max(r(s_2, \text{"stay"}) + \gamma V_2(s_2), r(s_2, \text{"go"}))] = [\max(-1 + 1\gamma, -2 + 3\gamma), \max(-1 + 3\gamma, 3)] = [\max(0, 1), \max(2, 3)] = [1, 3]$$

The optimal V is V_2 or V_3 because they both provide the highest value returns for each state across all iterations of V . With that being said, V_3 can generally be seen as better, since we show that the values have converged, whereas if we stopped at V_2 , we don't have any idea if the values were already their optimal values or not.

2 Value Iteration Convergence

$$1. \|V^0 - V^*\|_\infty = \|[-1, -3]\|_\infty = \max(|-1|, |-3|) = \max(1, 3) = 3$$

$$\|V^1 - V^*\|_\infty = \|[-2, 0]\|_\infty = \max(|-2|, |0|) = \max(2, 0) = 2$$

$$\|V^2 - V^*\|_\infty = \|[0, 0]\|_\infty = \max(|0|, |0|) = \max(0, 0) = 0$$

$$\|V^3 - V^*\|_\infty = \|[0, 0]\|_\infty = \max(|0|, |0|) = \max(0, 0) = 0$$

Clearly, the error decreases monotonically.

$$\begin{aligned}
2. \quad & \|T(V) - T(V')\|_\infty \\
&= \left\| \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V(s')] - \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V'(s')] \right\|_\infty \\
&\leq \max_a \left\| \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V(s')] - \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V'(s')] \right\|_\infty \\
&= \max_a \gamma \left\| \sum_{s'} p(s'|s, a) [r(s, a) + V(s')] - \sum_{s'} p(s'|s, a) [r(s, a) + V'(s')] \right\|_\infty \\
&= \max_a \gamma \sum_{s'} p(s'|s, a) \left\| [r(s, a) + V(s')] - [r(s, a) + V'(s')] \right\|_\infty \\
&= \max_a \gamma \sum_{s'} p(s'|s, a) \|V(s') - V'(s')\|_\infty \\
&\leq \max_a \gamma \sum_{s'} p(s'|s, a) \|V - V'\|_\infty \\
&= \gamma \|V - V'\|_\infty \max_a \sum_{s'} p(s'|s, a) \\
&= \gamma \|V - V'\|_\infty \quad (\text{We know } \sum_{s'} p(s'|s, a) = 1)
\end{aligned}$$

\therefore We have shown $\|T(V) - T(V')\|_\infty \leq \gamma \|V - V'\|_\infty$ \square

3. We want to show $\forall \epsilon > 0, \exists N > 0 \text{ s.t. } \forall n > N \quad \|V^{n+1} - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \epsilon$. This is shown below:

$$\begin{aligned}
& \|V^{n+1} - V^*\|_\infty \\
&= \|T(V^n) - T(V^*)\|_\infty \\
&\leq \gamma \|V^n - V^*\|_\infty \quad (\text{from the proof in question 2.2}) \\
&\leq \frac{\gamma}{1-\gamma} \|V^n - V^*\|_\infty \quad (\text{we assume } \gamma \text{ is between } 0 \text{ and } 1)
\end{aligned}$$

From the above, we've shown $\|V^{n+1} - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|V^n - V^*\|_\infty$. This indicates the distance between V^{n+1} and V^* shrinks over time (i.e. converges). As such, $\exists n$ s.t. $\|V^n - V^*\|_\infty \leq \epsilon$ (this n , in practice, is usually found when your program stops, indicating $\|V^n - V^{n-1}\|_\infty \leq \epsilon$). With this property, we see the following:

$$\begin{aligned}
& \|V^{n+1} - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|V^n - V^*\|_\infty \\
&\implies \|V^{n+1} - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \epsilon
\end{aligned}$$

\therefore We have shown $\|V^{n+1} - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \epsilon \quad \square$

4. Did not do this bonus question.

3 Learning the Model

1. Did not do this bonus question.

2. Did not do this bonus question.

3. Did not do this bonus question.

4. Did not do this bonus question.

4 Policy Gradients Variance Reduction

1. Let the approximation of the policy gradient be $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log \pi_{\theta}(\tau_i)$.

Now, let's show when $R(\tau) := R(\tau) - b$ does not change this estimate:

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log \pi_{\theta}(\tau_i) \\ \implies & \frac{1}{N} \sum_{i=1}^N (R(\tau_i) - b) \nabla_{\theta} \log \pi_{\theta}(\tau_i) \\ = & \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) \right] \left[\sum_{i=1}^N R(\tau_i) - b \right] \\ = & \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) R(\tau_i) \right] - \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) b \right] \end{aligned}$$

Now, in order to prove $\frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log \pi_{\theta}(\tau_i) = \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) R(\tau_i) \right] - \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) b \right]$, we must show $\frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) b = 0$. This can be seen below:

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) b \\ = & \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\nabla_{\theta} \log \pi_{\theta}(\tau_i) b \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[\mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} \left[\nabla_{\theta} \log \pi_{\theta}(\tau_i) b \right] \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot \mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} \left[\nabla_{\theta} \log \pi_{\theta}(\tau_i) \right] \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot \mathbb{E}_{a_t} \left[\nabla_{\theta} \log \pi_{\theta}(\tau_i) \right] \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot \int \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \pi_{\theta}(a_t | s_t) da_t \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot \nabla_{\theta} \int \pi_{\theta}(a_t | s_t) da_t \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot \nabla_{\theta} 1 \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[b \cdot 0 \right] \\ = & \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[0 \right] \\ = & 0 \end{aligned}$$

In the above mini-proof, we have shown for any t , the product of the gradient with b is 0.

As such, since the second term of $\frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) R(\tau_i) \right] - \frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) b \right]$ is 0, we can reduce it to $\frac{1}{N} \left[\sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) R(\tau_i) \right]$, and we observe that $\frac{1}{N} \sum_{i=1}^N (R(\tau_i) - b) \nabla_{\theta} \log \pi_{\theta}(\tau_i) = \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log \pi_{\theta}(\tau_i)$. \square

2. We are first going to calculate the variance $\text{Var}(\frac{1}{N} \sum_{i=1}^N (R(\tau_i) - b) \nabla_{\theta} \log \pi_{\theta}(\tau_i))$. We can use the rule that $\text{Var}(x) = \text{E}[x^2] - \text{E}[x]^2$ to solve this:

$$\begin{aligned}
& \text{Var}(\frac{1}{N} \sum_{i=1}^N (R(\tau_i) - b) \nabla_{\theta} \log \pi_{\theta}(\tau_i)) \\
&= \text{Var}(\text{E}_{\tau \sim \pi_{\theta}(\tau)} [(R(\tau) - b) \nabla_{\theta} \log \pi_{\theta}(\tau)]) \\
&= \text{Var}(\text{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b)]) \\
&= \text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] - \text{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b)]^2 \\
&= \text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] - \text{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) R(\tau)]^2 \quad (\text{Baseline is unbiased, as we showed in the previous question}) \\
&= \text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] - \text{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)]^2
\end{aligned}$$

As such, the variance is $\text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] - \text{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)]^2$. We can see the baseline will impact the first term's values, reducing its values and thus reducing the variance. As such, subtracting b helps reduce the variance of $\nabla_{\theta} J(\theta)$.

Now, we will calculate the baseline value leading to the least variance. To do this, we need to calculate the gradient of the variance with respect to the baseline b and set it to 0 and solve for b :

$$\begin{aligned}
& \frac{\delta \text{Var}}{\delta b} \\
&= \frac{\delta}{\delta b} \text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] - \text{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)]^2 \\
&= \frac{\delta}{\delta b} \text{E}_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] \quad (\text{second term doesn't depend on } b) \\
&= \frac{\delta}{\delta b} \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau) (R(\tau) - b))^2] \\
&= \frac{\delta}{\delta b} \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 (R(\tau) - b)^2] \\
&= \frac{\delta}{\delta b} \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 (R(\tau)^2 - 2R(\tau)b + b^2)] \\
&= \frac{\delta}{\delta b} \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)^2] - 2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 (R(\tau)b)] + \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 (b^2)] \\
&= \frac{\delta}{\delta b} \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)^2] - 2b\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] + b^2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] \\
&= \frac{\delta}{\delta b} (-2b)\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] + b^2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] \quad (\text{first term doesn't depend on } b) \\
&= -2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] + 2b\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] = 0 \\
&\implies -2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] + 2b\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] = 0 \\
&\implies 2b\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] = 2\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] \\
&\implies b\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2] = \text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)] \\
&\implies b = \frac{\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)]}{\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2]}
\end{aligned}$$

As such, the value of b that reduces the variance the most is $b = \frac{\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 R(\tau)]}{\text{E}[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2]}$.

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.

```
In [2]: %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`

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!

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

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.

```
In [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` .

```
In [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
```

[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 π^* . We can alternatively compute the value function V^π corresponding to an arbitrary policy π , 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^π , 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^π :

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

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^π 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.

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

```
In [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
```


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.

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

```
In [132]: 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
```

[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.

```
In [133]: #####
# Use this space for debugging #
# Make sure to delete this code before submission #
#####
pass
#####
```

```
In [134]: 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
```

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

Deliverable 1 (10 points)

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

```
In [135]: 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%
```

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.

```
In [136]: 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:       0.736
> Percentage of episodes goal reached:     93%
```

Deliverable 3 (5 bonus points)

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

```
In [137]: 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%
```

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.

```
In [138]: 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: 0.727

> Percentage of episodes goal reached: 93%

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)

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

```
In [1]: %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')
```

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 ϵ -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

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

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 [211]: from configs.p1_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: -1.00 +/- 0.00

1001/10000 [==>.....] - ETA: 13s - Loss: 3.8926 - Avg_R: 0.6650 - Max_R: 2.0000 - eps: 0.8020 - Grads: 12.9740 - Max_Q: 0.8200 - lr: 0.0042 - ETA: 11s - Loss: 1.0115 - Avg_R: 0.4400 - Max_R: 2.1000 - eps: 0.9208 - Grads: 5.3043 - Max_Q: 0.4004 - lr: 0 - ETA: 13s - Loss: 10.8250 - Avg_R: 0.3350 - Max_R: 3.0000 - eps: 0.8614 - Grads: 22.6048 - Max_Q: 0.6580 - lr: 0

Evaluating...

Average reward: 3.80 +/- 0.00

2001/10000 [=====>.....] - ETA: 13s - Loss: 12.6671 - Avg_R: 1.4500 - Max_R: 4.0000 - eps: 0.6238 - Grads: 32.2130 - Max_Q: 1.7724 - lr: 0.003 - ETA: 13s - Loss: 12.1808 - Avg_R: 1.1550 - Max_R: 3.8000 - eps: 0.6040 - Grads: 24.8829 - Max_Q: 1.8560 - lr: 0.0034

Evaluating...

Average reward: 3.90 +/- 0.00

3001/10000 [=====>.....] - ETA: 11s - Loss: 13.5121 - Avg_R: 1.9500 - Max_R: 4.1000 - eps: 0.4060 - Grads: 26.9340 - Max_Q: 2.5087 - lr: 0.0026 - ETA: 13s - Loss: 11.4555 - Avg_R: 1.7400 - Max_R: 3.8000 - eps: 0.5644 - Grads: 25.4581 - Max_Q: 2.0379 - lr: 0.00 - ETA: 12s - Loss: 15.0758 - Avg_R: 1.5700 - Max_R: 4.1000 - eps: 0.5248 - Grads: 29.2509 - Max_Q: 2.1887 - lr: 0.003 - ETA: 12s - Loss: 7.5629 - Avg_R: 1.4200 - Max_R: 4.0000 - eps: 0.5050 - Grads: 23.9891 - Max_Q: 2.2549 - lr:

Evaluating...

Average reward: 3.90 +/- 0.00

4001/10000 [=====>.....] - ETA: 10s - Loss: 16.1194 - Avg_R: 2.9150 - Max_R: 4.1000 - eps: 0.2080 - Grads: 25.7422 - Max_Q: 2.5906 - lr: 0.0018

Evaluating...

Average reward: 4.10 +/- 0.00

5001/10000 [=====>.....] - ETA: 8s - Loss: 3.4609 - Avg_R: 3.9850 - Max_R: 4.1000 - eps: 0.0100 - Grads: 12.1087 - Max_Q: 2.6064 - lr: 0.0010 - ETA: 10s - Loss: 9.0879 - Avg_R: 3.1500 - Max_R: 4.1000 - eps: 0.1684 - Grads: 23.9739 - Max_Q: 2.5935 - lr: - ETA: 9s - Loss: 2.3762 - Avg_R: 3.8700 - Max_R: 4.1000 - eps: 0.0892 - Grads: 13.5562 - Max_Q: 2.6092 - lr: 0.

Evaluating...

Average reward: 4.10 +/- 0.00

6001/10000 [=====>.....] - ETA: 6s - Loss: 0.2180 - Avg_R: 4.0500 - Max_R: 4.1000 - eps: 0.0100 - Grads: 4.8482 - Max_Q: 2.9011 - lr: 0.0010 - ETA: 7s - Loss: 3.0691 - Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 9.2268 - Max_Q: 2.7857 - lr

Evaluating...

Average reward: 4.10 +/- 0.00


```
7001/10000 [=====>.....] - ETA: 5s - Loss: 1.6149 - Avg_
R: 4.0950 - Max_R: 4.1000 - eps: 0.0100 - Grads: 10.1882 - Max_Q: 2.6138 - l
r: 0.001 - ETA: 5s - Loss: 0.4080 - Avg_R: 4.0100 - Max_R: 4.1000 - eps: 0.01
00 - Grads: 11.1709 - Max_Q: 2.6117 - lr: 0.0010
```

Evaluating...

Average reward: 3.80 +/- 0.00

```
8001/10000 [=====>.....] - ETA: 3s - Loss: 0.0142 - Avg_
R: 4.0050 - Max_R: 4.1000 - eps: 0.0100 - Grads: 2.3713 - Max_Q: 2.7571 - lr:
0.0010 - ETA: 4s - Loss: 0.2103 - Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100
- Grads: 6.0000 - Max_Q: 2.6040 - l
```

Evaluating...

Average reward: 4.10 +/- 0.00

```
9001/10000 [=====>...] - ETA: 1s - Loss: 0.2791 - Avg_
R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 3.5885 - Max_Q: 2.7440 - lr:
0.0010
```

Evaluating...

Average reward: 4.10 +/- 0.00

```
10001/10000 [=====] - ETA: 0s - Loss: 0.3712 - Avg_
R: 4.0950 - Max_R: 4.1000 - eps: 0.0100 - Grads: 9.1129 - Max_Q: 2.5361 - lr:
0.00 - 17s - Loss: 0.0600 - Avg_R: 4.0350 - Max_R: 4.1000 - eps: 0.0100 - Gra
ds: 1.8517 - Max_Q: 2.5327 - 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>)

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 [231]: 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: -0.10 +/- 0.00

Populating the memory 150/200...

Evaluating...

Average reward: -0.50 +/- 0.00

301/1000 [=====>.....] - ETA: 2s - Loss: 1.3270 - Avg_R: 0.1000 - Max_R: 2.8000 - eps: 0.4060 - Grads: 17.3416 - Max_Q: 0.0377 - lr: 0.0002

Evaluating...

Average reward: -0.50 +/- 0.00

401/1000 [=====>.....] - ETA: 2s - Loss: 2.2181 - Avg_R: 0.5150 - Max_R: 2.0000 - eps: 0.2080 - Grads: 42.4171 - Max_Q: 0.0895 - lr: 0.0001

Evaluating...

Average reward: 0.50 +/- 0.00

501/1000 [=====>.....] - ETA: 2s - Loss: 1.2275 - Avg_R: 0.5400 - Max_R: 2.0000 - eps: 0.0100 - Grads: 48.6401 - Max_Q: 0.1099 - lr: 0.0001

Evaluating...

Average reward: 0.50 +/- 0.00

601/1000 [=====>.....] - ETA: 2s - Loss: 5.6025 - Avg_R: 0.4500 - Max_R: 0.5000 - eps: 0.0100 - Grads: 47.5985 - Max_Q: 0.1239 - lr: 0.0001

Evaluating...

Average reward: 4.00 +/- 0.00

701/1000 [=====>.....] - ETA: 1s - Loss: 3.1406 - Avg_R: 3.8050 - Max_R: 4.0000 - eps: 0.0100 - Grads: 30.1634 - Max_Q: 0.2419 - lr: 0.0001

Evaluating...

Average reward: 3.90 +/- 0.00

801/1000 [=====>.....] - ETA: 1s - Loss: 0.9577 - Avg_R: 3.8450 - Max_R: 4.1000 - eps: 0.0100 - Grads: 20.7538 - Max_Q: 0.3436 - lr: 0.0001

Evaluating...

Average reward: 4.10 +/- 0.00

```

901/1000 [=====>...] - ETA: 0s - Loss: 2.5619 - Avg_R:
3.6300 - Max_R: 4.1000 - eps: 0.0100 - Grads: 53.5469 - Max_Q: 0.4259 - lr:
0.0001

Evaluating...
Average reward: 4.00 +/- 0.00

1001/1000 [=====] - 5s - Loss: 1.3149 - Avg_R: 3.850
0 - Max_R: 4.1000 - eps: 0.0100 - Grads: 144.1170 - Max_Q: 0.4986 - lr: 0.000
1

- 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> (<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 [234]: 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...

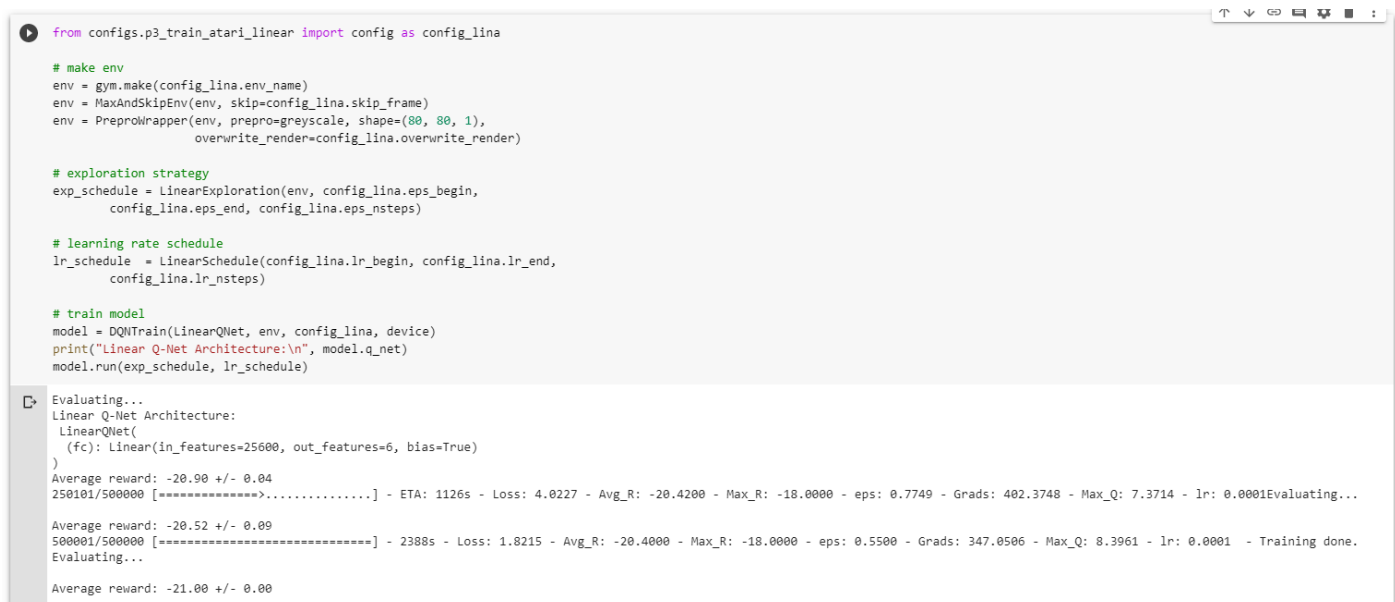
Linear Q-Net Architecture: LinearQNet((fc): Linear(in_features=25600, out_features=6, bias=True))

Average reward: -20.90 +/- 0.04 250101/500000 [=====>.....] - ETA: 1126s - Loss: 4.0227 - Avg_R: -20.4200 - Max_R: -18.0000 - eps: 0.7749 - Grads: 402.3748 - Max_Q: 7.3714 - lr: 0.0001Evaluating...

Average reward: -20.52 +/- 0.09 500001/500000 [=====] - 2388s - Loss: 1.8215 - Avg_R: -20.4000 - Max_R: -18.0000 - eps: 0.5500 - Grads: 347.0506 - Max_Q: 8.3961 - lr: 0.0001 - Training done. Evaluating...

Average reward: -21.00 +/- 0.00

NOTE: To make this code run faster, I used Google Colab. Above is the raw text I copy-pasted from the Colab notebook. Below is a screenshot of the results, which contains the same information. I just didn't want to waste my time running it all on my local machine.



```

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(
 (fc): Linear(in_features=25600, out_features=6, bias=True)
)
 Average reward: -20.90 +/- 0.04
 250101/500000 [=====>.....] - ETA: 1126s - Loss: 4.0227 - Avg_R: -20.4200 - Max_R: -18.0000 - eps: 0.7749 - Grads: 402.3748 - Max_Q: 7.3714 - lr: 0.0001Evaluating...
 Average reward: -20.52 +/- 0.09
 500001/500000 [=====] - 2388s - Loss: 1.8215 - Avg_R: -20.4000 - Max_R: -18.0000 - eps: 0.5500 - Grads: 347.0506 - Max_Q: 8.3961 - lr: 0.0001 - Training done.
 Evaluating...
 Average reward: -21.00 +/- 0.00

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 [ ]: 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)
```

NOTE: To make this code run faster, I used Google Colab. Below is a screenshot of the results. I just didn't want to waste my time running it all on my local machine. You can see the max reward achieved is -14.0

```
print("CNN Q-Net Architecture", model_q_net)
model.run(exp_schedule, lr_schedule)

(1): ReLU(
  (2): Conv2d(32, 64, kernel_size=(4, 4), stride=(2, 2))
  (3): ReLU(
  (4): Conv2d(64, 64, kernel_size=(4, 4), stride=(1, 1))
  (fc): Linear(in_features=1600, out_features=4, bias=True)
)
Average reward: -20.84 +/- 0.05
250001/500000 [>.....] - ETA: 21708s - Loss: 0.3046 - Avg_R: -20.4400 - Max_R: -19.0000 - eps: 0.7750 - Grads: 9.4592 - Max_Q: -0.3365 - lr: 0.0002 Evaluating...
Average reward: -20.34 +/- 0.10
500101/500000 [==>.....] - ETA: 21663s - Loss: 0.1064 - Avg_R: -20.1400 - Max_R: -17.0000 - eps: 0.5499 - Grads: 5.6625 - Max_Q: -0.2022 - lr: 0.0002 Evaluating...
Average reward: -19.90 +/- 0.13
750101/500000 [====>.....] - ETA: 20848s - Loss: 0.2746 - Avg_R: -19.3400 - Max_R: -16.0000 - eps: 0.3240 - Grads: 12.5154 - Max_Q: -0.2994 - lr: 0.0002 Evaluating...
Average reward: -19.66 +/- 0.16
1000101/500000 [=====>.....] - ETA: 19847s - Loss: 0.1616 - Avg_R: -19.3600 - Max_R: -16.0000 - eps: 0.1000 - Grads: 9.4356 - Max_Q: -0.2788 - lr: 0.0002 Evaluating...
Average reward: -18.52 +/- 0.20
1250001/500000 [=====>.....] - ETA: 18826s - Loss: 0.2194 - Avg_R: -19.0600 - Max_R: -17.0000 - eps: 0.1000 - Grads: 15.8378 - Max_Q: -0.2694 - lr: 0.0001 Evaluating...
Average reward: -20.00 +/- 0.12
1501101/500000 [=====>.....] - ETA: 17686s - Loss: 0.6278 - Avg_R: -18.3400 - Max_R: -15.0000 - eps: 0.1000 - Grads: 30.4464 - Max_Q: -0.2814 - lr: 0.0001 Evaluating...
Average reward: -19.48 +/- 0.16
1751401/500000 [=====>.....] - ETA: 16557s - Loss: 0.3256 - Avg_R: -19.0400 - Max_R: -14.0000 - eps: 0.1000 - Grads: 20.5341 - Max_Q: -0.2397 - lr: 0.0001 Evaluating...
Average reward: -18.78 +/- 0.20
2001701/500000 [=====>.....] - ETA: 15419s - Loss: 0.4824 - Avg_R: -18.7200 - Max_R: -14.0000 - eps: 0.1000 - Grads: 28.7398 - Max_Q: -0.2451 - lr: 0.0001 Evaluating...
Average reward: -18.54 +/- 0.18
2252001/500000 [=====>.....] - ETA: 14212s - Loss: 0.2669 - Avg_R: -18.6400 - Max_R: -15.0000 - eps: 0.1000 - Grads: 18.9994 - Max_Q: -0.2789 - lr: 0.0000 Evaluating...
Average reward: -18.74 +/- 0.19
2502301/500000 [=====>.....] - ETA: 12973s - Loss: 0.4094 - Avg_R: -18.7400 - Max_R: -15.0000 - eps: 0.1000 - Grads: 19.7147 - Max_Q: -0.3259 - lr: 0.0000 Evaluating...
Average reward: -18.60 +/- 0.26
2752701/500000 [=====>.....] - ETA: 11726s - Loss: 0.1200 - Avg_R: -18.6800 - Max_R: -16.0000 - eps: 0.1000 - Grads: 9.1460 - Max_Q: -0.3171 - lr: 0.0000 Evaluating...
Average reward: -18.18 +/- 0.22
3003001/500000 [=====>.....] - ETA: 10470s - Loss: 0.2166 - Avg_R: -18.6800 - Max_R: -15.0000 - eps: 0.1000 - Grads: 18.9064 - Max_Q: -0.3340 - lr: 0.0000 Evaluating...
Average reward: -18.36 +/- 0.25
3253401/500000 [=====>.....] - ETA: 9191s - Loss: 1.0473 - Avg_R: -18.5800 - Max_R: -15.0000 - eps: 0.1000 - Grads: 30.2055 - Max_Q: -0.3231 - lr: 0.0000 Evaluating...
Average reward: -18.82 +/- 0.19
3503501/500000 [=====>.....] - ETA: 7900s - Loss: 0.6123 - Avg_R: -18.8400 - Max_R: -14.0000 - eps: 0.1000 - Grads: 27.9635 - Max_Q: -0.3253 - lr: 0.0000 Evaluating...
Average reward: -18.28 +/- 0.22
3753701/500000 [=====>.....] - ETA: 6579s - Loss: 0.1775 - Avg_R: -18.2800 - Max_R: -14.0000 - eps: 0.1000 - Grads: 13.6898 - Max_Q: -0.3307 - lr: 0.0000 Evaluating...
Average reward: -18.56 +/- 0.22
4004001/500000 [=====>.....] - ETA: 5258s - Loss: 0.1882 - Avg_R: -18.7800 - Max_R: -15.0000 - eps: 0.1000 - Grads: 16.4965 - Max_Q: -0.3412 - lr: 0.0000 Evaluating...
Average reward: -18.64 +/- 0.23
4254301/500000 [=====>.....] - ETA: 3944s - Loss: 0.0999 - Avg_R: -18.2600 - Max_R: -14.0000 - eps: 0.1000 - Grads: 13.6329 - Max_Q: -0.3464 - lr: 0.0000 Evaluating...
Average reward: -18.20 +/- 0.25
4454501/500000 [=====>.....] - ETA: 2884s - Loss: 0.1803 - Avg_R: -18.3800 - Max_R: -14.0000 - eps: 0.1000 - Grads: 20.3310 - Max_Q: -0.3332 - lr: 0.0000 Buffered data was truncated after reaching the output size limit.
```

In []: