

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 [3]: %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 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 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

$$Q_{sample}(s_t) = r_t \text{ if done} \\ = r_t + \gamma \max_{a_{t+1}} Q_{target}(s_{t+1}, a_{t+1}) \text{ otherwise}$$

$$\text{Loss} = (Q_{sample}(s_t) - Q(s_t, a_t))^2$$

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 [4]: 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: 1.60 +/- 0.00

1001/10000 [==>.....] - ETA: 3s - Loss: 0.2968 - Avg_R: 0.5500 - Max_R: 3.1000 - eps: 0.8020 - Grads: 1.2544 - Max_Q: 0.6897
- lr: 0.0042

Evaluating...
Average reward: 3.80 +/- 0.00

2001/10000 [=====>.....] - ETA: 3s - Loss: 0.3410 - Avg_R: 2.0550 - Max_R: 4.1000 - eps: 0.6040 - Grads: 1.1430 - Max_Q: 1.8823
- lr: 0.0034

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

3001/10000 [=====>.....] - ETA: 3s - Loss: 0.3964 - Avg_R: 2.0800 - Max_R: 4.0000 - eps: 0.4060 - Grads: 0.7165 - Max_Q: 2.5491
- lr: 0.0026

Evaluating...
Average reward: 3.80 +/- 0.00

4001/10000 [=====>.....] - ETA: 2s - Loss: 0.3749 - Avg_R: 2.9400 - Max_R: 4.1000 - eps: 0.2080 - Grads: 0.6549 - Max_Q: 2.7646
- lr: 0.0018

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

5001/10000 [=====>.....] - ETA: 2s - Loss: 0.0746 - Avg_R: 3.9900 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.5588 - Max_Q: 2.7375
- lr: 0.0010

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

6001/10000 [=====>.....] - ETA: 1s - Loss: 0.0002 - Avg_R: 3.9000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.0812 - Max_Q: 2.9415
- lr: 0.0010

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

7001/10000 [=====>.....] - ETA: 1s - Loss: 0.0001 - Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.0516 - Max_Q: 2.6616
- lr: 0.0010

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

8001/10000 [=====>.....] - ETA: 0s - Loss: 0.0055 - Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.1861 - Max_Q: 2.8178
- lr: 0.0010

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

9001/10000 [=====>...] - ETA: 0s - Loss: 0.0001 - Avg_R: 4.0000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.0183 - Max_Q: 2.7897
- lr: 0.0010

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

10001/10000 [=====] - 4s - Loss: 0.0001 - Avg_R: 4.1000 - Max_R: 4.1000 - eps: 0.0100 - Grads: 0.0306 - Max_Q: 2.8013 - 1
r: 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 [20]: 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: 0.50 +/- 0.00

Populating the memory 150/200...

Evaluating...
Average reward: 0.50 +/- 0.00

 301/1000 [=====>.....] - ETA: 2s - Loss: 0.0899 - Avg_R: 0.3250 - Max_R: 3.1000 - eps: 0.4060 - Grads: 4.2915 - Max_Q: 0.1761 -
lr: 0.0002

Evaluating...
Average reward: 0.50 +/- 0.00

 401/1000 [=====>.....] - ETA: 2s - Loss: 0.0347 - Avg_R: 0.0750 - Max_R: 2.3000 - eps: 0.2080 - Grads: 1.3289 - Max_Q: 0.2084 -
lr: 0.0001

Evaluating...
Average reward: 0.50 +/- 0.00

 501/1000 [=====>.....] - ETA: 2s - Loss: 0.0188 - Avg_R: 0.4900 - Max_R: 1.9000 - eps: 0.0100 - Grads: 1.0697 - Max_Q: 0.2152 -
lr: 0.0001

Evaluating...
Average reward: 0.50 +/- 0.00

 601/1000 [=====>.....] - ETA: 2s - Loss: 0.2072 - Avg_R: 2.9450 - Max_R: 4.0000 - eps: 0.0100 - Grads: 3.6084 - Max_Q: 0.3536 -
lr: 0.0001

Evaluating...
Average reward: 3.80 +/- 0.00

 701/1000 [=====>.....] - ETA: 1s - Loss: 0.0938 - Avg_R: 3.9250 - Max_R: 4.0000 - eps: 0.0100 - Grads: 3.8625 - Max_Q: 0.4915 -
lr: 0.0001

Evaluating...
Average reward: 3.80 +/- 0.00

 801/1000 [=====>.....] - ETA: 1s - Loss: 0.0779 - Avg_R: 3.8500 - Max_R: 4.0000 - eps: 0.0100 - Grads: 1.8714 - Max_Q: 0.6081 -
lr: 0.0001

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

 901/1000 [=====>...] - ETA: 0s - Loss: 0.0088 - Avg_R: 3.9750 - Max_R: 4.0000 - eps: 0.0100 - Grads: 2.9653 - Max_Q: 0.6976 -
lr: 0.0001

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

1001/1000 [=====] - 7s - Loss: 0.0059 - Avg_R: 3.8250 - Max_R: 4.0000 - eps: 0.0100 - Grads: 1.9607 - Max_Q: 0.7689 - lr:
0.0001

- Training done.
Evaluating...
Average reward: 4.00 +/- 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/> (<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 hyperparamters 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 [ ]: 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(
  (layer): Linear(in_features=25600, out_features=6, bias=True)
)
```

Average reward: -20.92 +/- 0.04

130601/500000 [=====.....] - ETA: 1647s - Loss: 0.1784 - Avg_R: -20.6000 - Max_R: -18.0000 - eps: 0.8825 - Grads: 17.9362 - Max_Q: 6.9595 - lr: 0.0001

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_gnet.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 stpes:

1. EAR >= 0.0 : 4/4 points
2. EAR >= -5.0 : 3/4 points
3. EAR >= -10.0 : 3/4 points
4. EAR >= -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)
```

```
In [ ]:
```