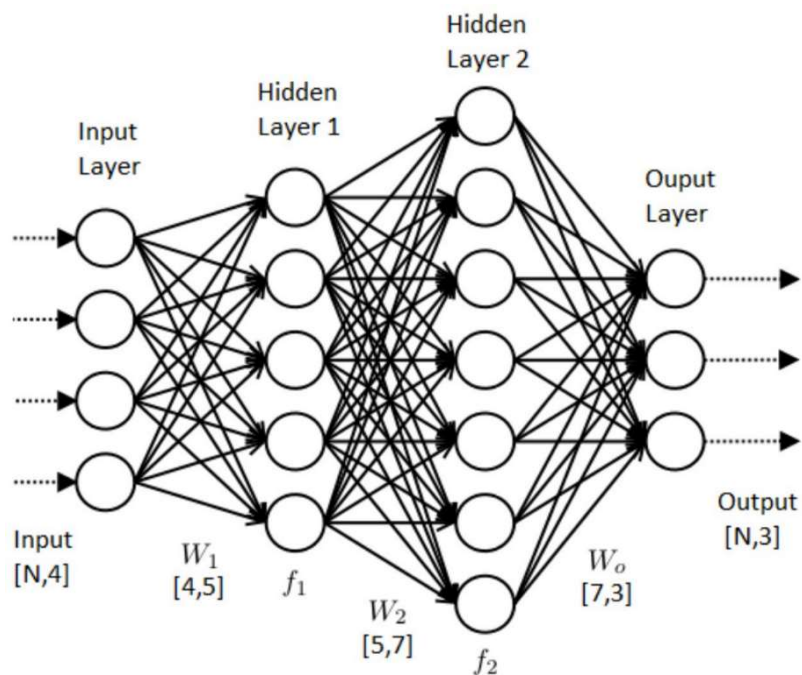# Solving Navigation using Deep Reinforcement Learning (Value-Based Methods)

By: Adrian Chow

Deep reinforcement learning is a subfield of machine learning that combines reinforcement learning and deep learning. Rather than using traditional RL methodologies, such as a Q-Table, artificial neural networks allow us to expand an agent's ability to solve a complex environment. While this post focuses on tackling the navigation issue, there is a lot to learn about how a DQN (Deep Q-Network) can use non-linear approximators to estimate the optimal action-value function.
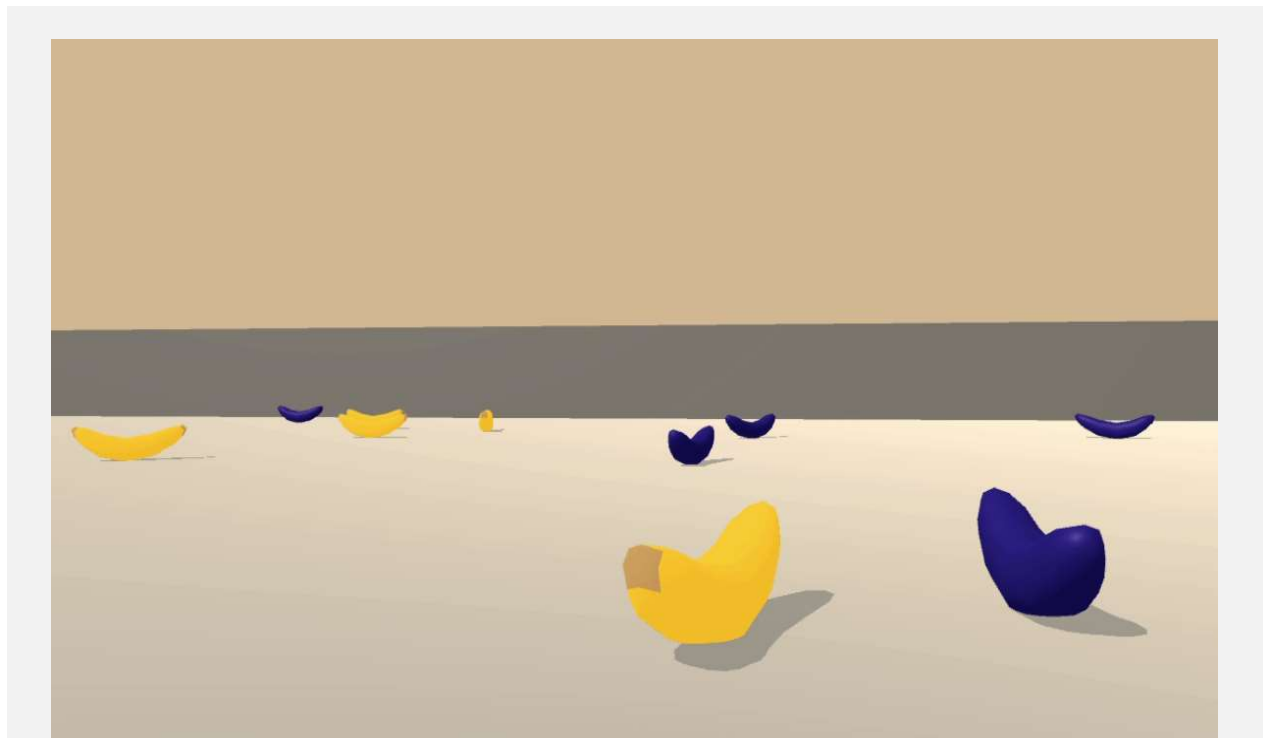


Artificial Neural Network (Source: VIASAT)

## The Environment

To solve the task of navigation, I have used a task similar to the Food Collector environment on the [Unity ML-Agents GitHub page](). Instead, the project is provided by Udacity as a part of their Deep Reinforcement Learning Nanodegree.

In this environment, an Agent navigates a large, square world collecting bananas. Each episode of this task is limited to 300 steps. A reward of +1 is provided for collecting a yellow banana, and a reward of -1 is provided for collecting a blue banana. Thus, the goal of the agent is to collect as many yellow bananas as possible, avoiding the blue ones.

RL Problem Specifications:

- The Goal of the Agent: Collect yellow bananas, avoid blue bananas

- Rewards : `+1` collecting yellow bananas, `-1` collecting blue bananas

- Action Space — Discrete, 4 Actions [Forward, Back, Left, Right]

- State Space — Continuous, 37-Dimension

- Solving Condition: Average score of +13 over 100 consecutive episodes.

The state-space has 37 dimensions and contains the agent's velocity, along with the ray-based perception of objects around the agent's forward direction. Given this information, the agent has to learn how to best select actions.
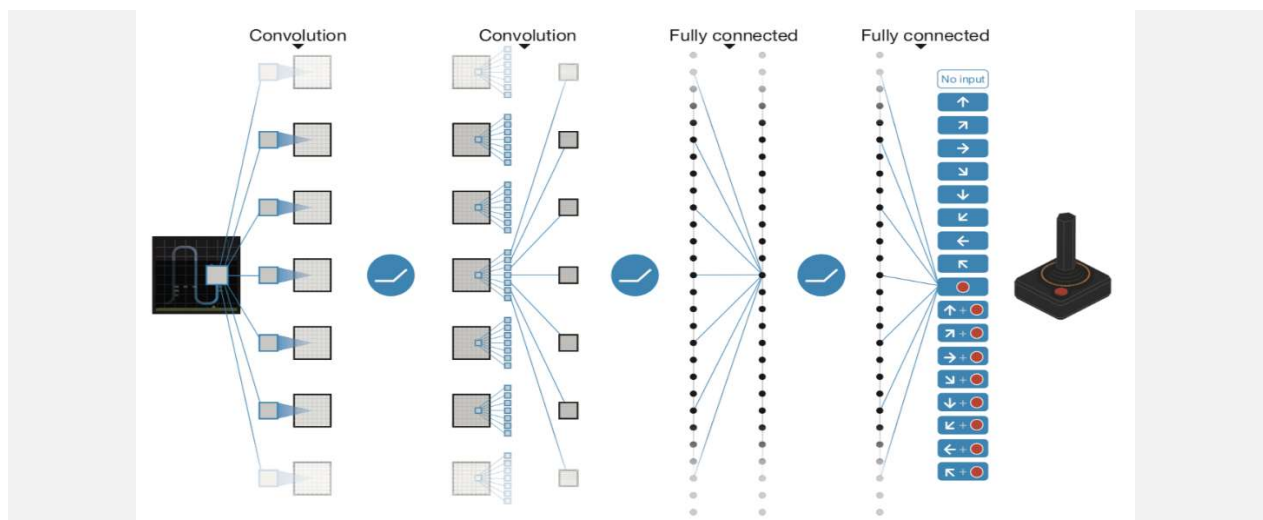
## Learning Algorithm



Illustration of DQN Architecture ([Source](#))

The learning algorithm was developed by researchers working towards achieving human-level control to play a challenging domain of classic Atari 2600 games! To learn more about this, you read this [research paper](#) which outlines the algorithm in full detail.

*Initialization:*

- Initialize Replay Buffer Data Structure. A memory bank/buffer with a capacity of N=10000 memories should suffice. This data structure is essential a deque that stores nodes (named tuples) with the following data: ["state", "action", "reward", "next_state", "done"]

- Initialize Local Model: A feed-forward NN that will estimate the action-value function.

- Initialize Target Model: A feed-forward NN that will estimate the target action-value function.

The implementation of the simple feed-forward model can be seen below. In addition, you can see the full implementation of the replay buffer here.

```python
1  import torch
2  import torch.nn as nn
3
4  class QNetwork(nn.Module):
5
6      def __init__(self, state_size, action_size, seed, fc1_units=64, fc2_units=64):
7
8          """
9
10         Actor (Policy) Model
11
12         Params
13         ======
14             state_size (int): Dimension of each state
15             action_size (int): Dimension of each action
16             seed (int): Random seed
17             fc1_units (int): Number of nodes in first hidden layer
18             fc2_units (int): Number of nodes in second hidden layer
19         """
20
21         # Inherit nn.Module as subclass
22         super(QNetwork, self).__init__()
23         self.seed = torch.manual_seed(seed)
24         self.action_size = action_size
25         self.state_size = state_size
26
27          # Define layer parameters
28         self.fully_connected_1 = nn.Linear(state_size, fc1_units)
29         self.fully_connected_2 = nn.Linear(fc1_units, fc2_units)
30         self.fully_connected_3 = nn.Linear(fc2_units, action_size)
31
32     def foward(self, state):
33         """
34         Foward Pass that maps state -> action values.
35
36         """
37         X = nn.functional.relu(self.fully_connected_1(state))
38         X = nn.functional.relu(self.fully_connected_2(X))
39         return self.fully_connected_3(X)
```

*For episode e ←1 to M:*

- Initialize the input frame x, a.k.a the initial state. Most OpenAI Gym environments will provide a simple API to grab a restated state.

```
for i_episode in range(1, n_episodes + 1):
    env_info = env.reset(train_mode=True)[brain_name]
    state = env_info.vector_observations[0]
```

*For step t ← 1 to T:*

- We break this down into two parts: **SAMPLE** and **LEARN**

```
1    for t in range(max_t):
2
3                # Use agent model to predict next action
4                action = agent.act(state, epsilon)
5
6                env_info = env.step(action)[brain_name]          # send the action to the environment
7                next_state = env_info.vector_observations[0]     # get the next state
8                reward = env_info.rewards[0]                      # get the reward
9                done = env_info.local_done[0]                     # see if episode has finished
10
11               agent.step(state, action, reward, next_state, done)
12
13               # Store results
14               state = next_state
15               score += reward
16
17               if done:
18                   break
```

forSteps.py hosted with ♥ by GitHub                                                          view raw

- Sample Step:

1. Choose action A from state S using policy π ← ε-Greedy( q(S, A, w) ).

2. Take Action A, observe reward R, and next input frame X{t+1}

3. Prepare next state and store experience tuple in Replay Buffer

4. Set S ← S'

```python
def step(self, state, action, reward, next_state, done):
    """

    Process a step from time step t to t+1 by updating agent models.

    Params
    ======
        state (continious): current state before action
        action (discrete): action take for given state
        reward (int): reward recieved after performing action
        next_state (int): state achieved at timestep t+1
        done (bool): episode completed on this timestep

    """

    self.memory.add(state, action, reward, next_state, done)

    # Increase counter until we are ready to take an update step
    self.t_step = (self.t_step + 1) % UPDATE_EVERY

    if self.t_step == 0:
        # If enough samples are available in memory, get random subset
        if len(self.memory) > BATCH_SIZE:
            experiences = self.memory.sample()
            self.learn(experiences, GAMMA)
```

- Learn Step:

1. Obtain random minibatch of tuples ("state", "action", "reward", "next_state", "done") from Replay Buffer

2. Set TD target and perform the update step on the local model

$$\Delta w = \alpha \cdot (\overbrace{R + \gamma \max_{a} \hat{q}(S', a, w^{-})}^{\text{TD error}} - \underbrace{\hat{q}(S, A, w)}_{\text{old value}}) \nabla_{w} \hat{q}(S, A, w)$$

$$\underbrace{\phantom{R + \gamma \max_{a} \hat{q}(S', a, w^{-})}}_{\text{TD target}}$$

Update step on weights

3. Perform a soft update step of the target model by factoring in a portion of the local model's parameters

```python
1    def learn(self, experiences, gamma):
2        """
3
4        Update value parameters using given batch of experience tuples.
5
6        Params
7        ======
8            experiences (Tuple[torch.Variable]): tuple of (s, a, r, s', done) tuples
9            gamma (float): discount factor
10
11        """
12
13        # Note: These tensors make up a batch of 64 experiences
14
15        states, actions, rewards, next_states, dones = experiences
16
17        # Get max predicted Q values (for next states) from target model
18        q_prime = self.qnetwork_target.foward(next_states)
19        # Choose the reward from action that gives max return
20        q_prime = q_prime.detach().max(1)[0].unsqueeze(1)
21
22        not_done_bool = (1 - dones) # If done, no need to include next return
23        td_target = rewards + (gamma * q_prime) * not_done_bool
24
25        # This is the model we will update
26        q_expected = self.qnetwork_local.foward(states)
27        # Gathers the expected values for each action
28        q_expected = q_expected.gather(1, actions)
29
30        # Compute the loss, minimize the loss
31        loss = F.mse_loss(td_target, q_expected)
32        self.optimizer.zero_grad() # reset gradient
33        loss.backward() # Calculate the gradient
34        self.optimizer.step()  # Update weights
35
36        # Update the target model parameters (Soft Update)
37        # Soft Update: Factor in local parameter changes by a factor of TAU
38        # Rather than update for every C steps, this helps inch closer to local parameters
39
40        for target_param, local_param in zip(self.qnetwork_target.parameters(),
41                                              self.qnetwork_local.parameters()):
42
43            upd_wghts = ((1.0 - self.tau) * target_param.data) + (self.tau * local_param.data)
44            target_param.data.copy_(upd_wghts)
```

Yeah, so that's it! The model is able to learn over time and develop the ability to collect yellow bananas while avoiding the blue bananas. The full repository can be viewed here.
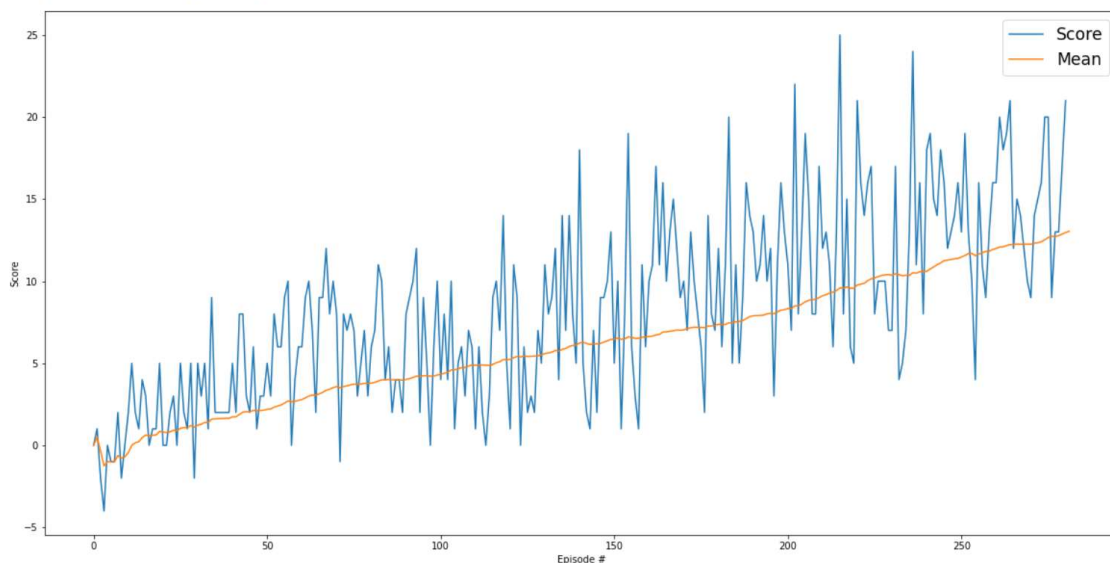
The application of this project is targeted for, while not limited to, robotic navigation. Some ideas that one could build off of this project are Garbage Collector Robots and Agricultural Robots (Muti-purpose)!

**Plot of Results**

Training Parameters:

- Max Episodes: 2000

- Max Time Steps: 1000

- Buffer Size: 10000

- Batch Size: 64

- Gamma: 0.99

- Tau: 1e-3

- Learning Rate = 5e-4

- Epsilon (Start, End, Decay Rate): [0.10, 0.01, 0.987]

```
Beginning Training....
Episode 100     Average Score: 4.28
Episode 200     Average Score: 8.26
Episode 282     Average Score: 13.04
Environment solved in 282 episodes!     Average Score: 13.04
eps: 0.01
last_score: 14.0
solved_in: 282
last_100_avg: 13.04
save_file: baseline_params.pth
```

**Going Further…**

Several improvements to the original Deep Q-Learning algorithm have been suggested.

1. **Double DQN**

Deep Q-Learning tends to overestimate action values. Double Q-Learning has been shown to work well in practice to help with this. The main idea is to gain a second opinion on the target action-value (TD Target) by running the state through the local model first, then the target model. See the changes below:

```
1          if self.double_DQN:
2
3              # Get max predicted Q values (for next states) from local model
4              q_prime = self.qnetwork_local.foward(next_states)
5              # For the batch, we want to store the most greedy action for Double DQN Networks
6              greedy_action_next = q_prime.max(dim=1,keepdim=True)[1]
7              # Choose the reward from action that gives max return
8              q_prime = q_prime.detach().max(1)[0].unsqueeze(1)
9
10             # For DDQN, we must run the perform a sanity check using both nets,
11             # while still taking the original greedy actions
12             DDQN_q_prime = self.qnetwork_target.foward(next_states)
13             DDQN_q_prime = DDQN_q_prime.gather(1, greedy_action_next)
14
15             not_done_bool = (1 - dones) # If done, no need to include next return
16             td_target = rewards + (gamma * DDQN_q_prime) * not_done_bool
```

DoubleDQN.py hosted with ♥ by GitHub                                    view raw

Double DQN Changes

## 2. Prioritized Experience Replay

Deep Q-Learning samples experience transitions *uniformly* from a replay memory. Prioritized experience replay is based on the idea that the agent can learn more effectively from some transitions than from others, and the more important transitions should be sampled with higher probability. By assigning priority we can steer the agent towards understanding the more important experiences.

```python
1   import numpy as np
2   import torch
3   import random
4   from collections import namedtuple, deque
5
6   version = 1.0
7   EXP_FIELD_NAME = ["state", "action", "reward", "next_state", "done", "prob"]
8
9   class PrioritizedReplayBuffer:
10      """
11      ReplayBuffer (Data Structure): Fixed-size buffer to store experience tuples.
12      Experience (Data): NamedTuple with structure of
13                          ["state", "action", "reward", "next_state", "done"]
14      """
15
16      def __init__(self,
17                  action_size,
18                  buffer_size,
19                  batch_size,
20                  seed,
21                  device,
22                  prioritized_params):
23          """
24          Initialize a ReplayBuffer object.
25
26          Params
27          ======
28              action_size (int): dimension of each action
29              buffer_size (int): maximum size of buffer
30              batch_size (int): size of each training batch
31              seed (int): random seed
32              device (string): compute resource
33              prioritized_params (dict): contains the constants for hyperparameters
34                α (int): α = 0 corresponding to the uniform case, α = 1 to using priorities
35                b (int): b is the bias used to scale the change in weights for importance-sampling 
36
37          """
38
39          self.action_size = action_size
40          self.buffer_size = buffer_size
41          self.batch_size = batch_size
```

```python
42            self.seed = seed
43            self.device = device
44            self.a = prioritized_params['a']
45            self.b = prioritized_params['b']
46            self.b_inc_rate = prioritized_params['b_inc_rate']
47            self.e = prioritized_params['e']
48            self.b_max = 1.0
49            self.max_prob = 1.0
50
51            # Memories will be stored on a double-ended queue (deque)
52            self.memory = deque(maxlen=self.buffer_size)
53            # Now we have to keep track of priorities as well
54            self.priorities = deque(maxlen=self.buffer_size)
55
56            # Experiences will be a named tupile structure
57            self.experience = namedtuple("Experience", field_names=EXP_FIELD_NAME)
58
59
60
61        def get_td_errors(self, state, action, reward, next_state, done, models):
62
63            """ Calculate td error."""
64            models['local'].eval() # Evaluation Mode
65            with torch.no_grad(): # No Gradient Descent
66
67                # Make 64 copies to fit into model
68                states = []
69                rewards = []
70                next_states = []
71                dones = []
72
73                for i in range(self.batch_size):
74                    states.append(state)
75                    rewards.append(reward)
76                    next_states.append(next_state)
77                    dones.append(done)
78
79                states = torch.from_numpy(np.vstack(states)).float().to(self.device)
80                q_expected = models['local'].foward(states)
81
82                next_states = torch.from_numpy(np.vstack(next_states)).float().to(self.device)
83                q_prime = models['target'].foward(next_states)
```

```python
85              rewards = torch.from_numpy(np.vstack(rewards)).float().to(self.device)
86              dones = torch.from_numpy(np.vstack(dones)).float().to(self.device)
87              not_done_bool = (1 - dones) # If done, no need to include next return
88              td_target =  rewards + (models['GAMMA'] * q_prime) * not_done_bool
89
90              td_error = td_target - q_prime
91
92          models['local'].train()
93
94          return abs(td_error.detach()[0][action])
95
96
97      def add(self, state, action, reward, next_state, done, models):
98          """
99
100         Append an experience to the memory.
101
102         Params
103         ======
104             state (continious): current state before action
105             action (discrete): action take for given state
106             reward (int): reward recieved after performing action
107             next_state (int): state achieved at timestep t+1
108             done (bool): episode completed on this timestep
109
110         """
111         # Create an experience
112
113         #Calculate the td_error
114         td_error = self.get_td_errors(state, action, reward, next_state, done, models)
115         td_error += self.e
116         exp = self.experience(state, action, reward, next_state, done, td_error)
117         self.memory.append(exp)
118         self.priorities.append(td_error)
119
120
121     def sample(self):
122         """
123
124         Randomly sample an experience from the memory.
125
126         return: experience (tuple)
127
128         """
```
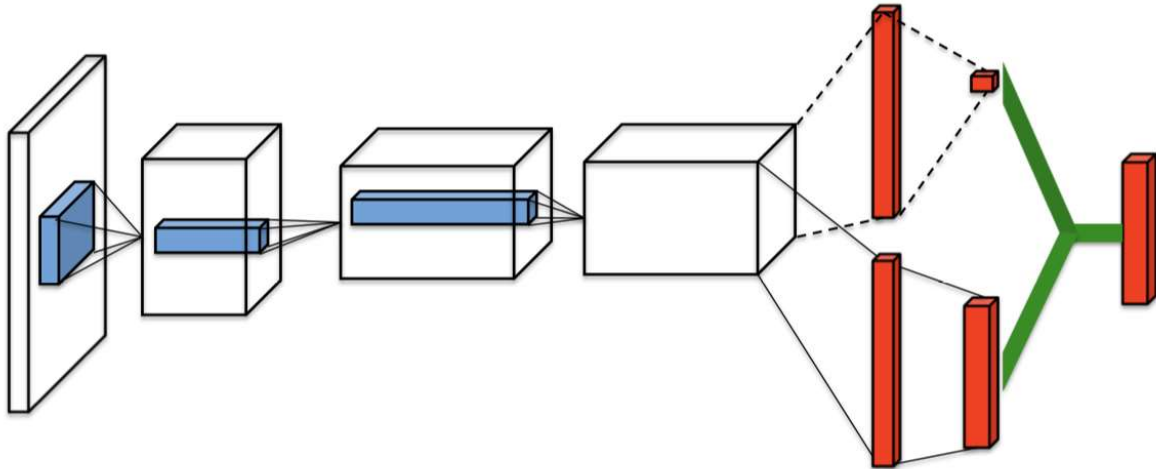
```python
130         #########################################
131         #   Prioritized Sampling Modifications   #
132         #########################################
133
134         # Calculate the Sampling priorities
135         sampling_prob = np.array(self.priorities)
136         sampling_prob = sampling_prob ** self.a / sum(sampling_prob ** self.a)
137         sample_idxs = np.random.choice(np.arange(len(self.memory)), size=self.batch_size, p=sampl
138         experiences = []
139         update_factors = []
140         for i in sample_idxs:
141             experiences.append(self.memory[i])
142             update_factors.append( ((1 / self.buffer_size ) * (1 / sampling_prob[i] ) )** self.b
143
144         # For each item in the tuple, stack vertically and convert to GPU torch tensor
145         states = np.vstack([e.state for e in experiences if e is not None])
146         states = torch.from_numpy(states).float().to(self.device) # (float)
147
148         actions = np.vstack([e.action for e in experiences if e is not None])
149         actions = torch.from_numpy(actions).long().to(self.device) # (long)
150
151         rewards = np.vstack([e.reward for e in experiences if e is not None])
152         rewards = torch.from_numpy(rewards).float().to(self.device) # (float)
153
154         next_states = np.vstack([e.next_state for e in experiences if e is not None])
155         next_states = torch.from_numpy(next_states).float().to(self.device) # float
156
157         dones = np.vstack([e.done for e in experiences if e is not None]).astype(np.uint8) # Make
158         dones = torch.from_numpy(dones).float().to(self.device)
159
160         update_factors = np.vstack(update_factors)
161         update_factors = torch.from_numpy(update_factors).float().to(self.device)
162
163         # Update b
164         self.b = min(self.b * self.b_inc_rate, self.b_max)
165
166         return (states, actions, rewards, next_states, dones, update_factors)
167
168     def __len__(self):
169         """
170
171         Return the current size of internal memory.
172
173         return: length (int)
174
175         """
```

PER Replay Buffer Implementation

As you can see the implementation is more involved than the previous replay buffer. It is important to note that prioritized experience replay adds additional time and space complexity to the algorithm.

**3. Dueling DQN**



Dueling DQN

Currently, in order to determine which states are (or are not) valuable, we have to estimate the corresponding action values *for each action*. However, by replacing the traditional Deep Q-Network (DQN) architecture with a dueling architecture, we can assess the value of each state, without having to learn the effect of each action. By using two streams (value and advantage), we can gain a better estimate of the action-value.

Here I have implemented my own version of the dueling network with the help of Chris Yoon's post, feel free to take a look.

```python
class DuelingQNetwork(nn.Module):

    def __init__(self, state_size, action_size, seed, fc1_units=64, fc2_units=64, fc3_units=32):


        """

        Actor (Policy) Model

        Params
        ======
            state_size (int): Dimension of each state
            action_size (int): Dimension of each action
            seed (int): Random seed
            fc1_units (int): Number of nodes in first hidden layer
            fc2_units (int): Number of nodes in second hidden layer
            fc2_units (int): Number of nodes in thirs hidden layer
        """

        # Inherit nn.Module as subclass
        super(DuelingQNetwork, self).__init__()
        self.seed = torch.manual_seed(seed)

        self.feauture_layer = nn.Sequential(
            nn.Linear(state_size, fc1_units),
            nn.ReLU(),
            nn.Linear(fc1_units, fc2_units),
            nn.ReLU()
        )

        self.value_stream = nn.Sequential(
            nn.Linear(fc2_units, fc3_units),
            nn.ReLU(),
            nn.Linear(fc3_units, action_size),
        )

        self.advantage_stream = nn.Sequential(
            nn.Linear(fc2_units, fc3_units),
            nn.ReLU(),
            nn.Linear(fc3_units, action_size),
        )


    def foward(self, state):
        """
        Foward Pass that maps state -> action values.

        """
        # Foward Pass
        features = self.feauture_layer(state)

        # Value Stream - State Values
        values = self.value_stream(features)
        # Advantage Stream - Advantage Values
        advantages  = self.advantage_stream(features)

        # Q Estimate = V(state) + (A(state) - mean(A(state))
        qvals = values + (advantages - advantages.mean())
```

# Ideas for Future Work

The next step for this project would be make the priority replay buffer more efficient. When training the algorithm, it seemed that the training time took too long. In reality, we want to have efficient algorithms in order to perform real-time computations. Maybe we could write the implementation in C++ or use a more efficient method of estimating priority.