

DQN Sort

Name : 藤原 大悟

Student Number : 6930-31-6255

Formulation

数字列の昇順ソートタスク

状態 s_t : N 桁の0~9の数字(x_1, x_2, \dots, x_N)の列(ただし0が先頭でも良いこととした)

行動 a_t : 数字の隙間($i=1 \sim N-1$)に対し, その両隣の数字を交換する操作. 確率1で交換後の一つの数字列を観測する.

報酬 r_t : 数列の昇順ソートに成功すると+1, それ以外の場合は一回の操作ごとに $-\frac{2}{N(N-1)}$ (罰則)

ニューロン数が N である隠れ層を二つ有する4層のニューラルネットワークでDQN強化学習を行う. 入力を状態 s_t , 出力を各行動 a_t の有する価値としてQ関数のニューラルネットワークによる近似を行う.

誤差はTD誤差

$$TDError = y_t - Q_{\theta}(s_t, a_t)$$

$$y_t = \begin{cases} r_t & s_{t+1} \text{が終端状態} \\ r_t + \max_{a'} \hat{Q}_{\theta}(s_{t+1}, a') & \text{otherwise} \end{cases}$$

を用い, これに対しHuberLossを考えた.

ターゲットネットワークの更新頻度Bを128とした.

ミニバッチサイズを16, キューのサイズを100としてExperience Replayを適用した.

最適化手法はAdamを採用した.

ϵ -greedy法のパラメータ ϵ は初期値を1とし, 以降指数的に各エポックごとに減衰し, 0に向かうようにした.

学習の開始は更新回数がA=1000を超えてからとした.

N=4とした

Implementation

Tensorflowを用いて実装を行なった

In [143]:

```
1 import matplotlib.pyplot as plt
2 from collections import deque
3 import tensorflow as tf
4 import numpy as np
5 import math
6 from scipy.stats import multivariate_normal as multi_gauss
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import warnings
11 warnings.simplefilter(action='ignore', category=FutureWarning)
```

In [147]:

```

1  class DQN_SORT:
2      def __init__(self,inputnum):
3          self.N=inputnum
4          #隠れ層1のニューロン数
5          self.M=100
6          #隠れ層2のニューロン数
7          self.L=100
8          #行動数
9          self.Nact=self.N-1
10         #ミニバッチ数
11         self.minibatch_size = 16
12         #キューのサイズ
13         self.replay_memory_size = 1000
14         self.epsilon=0.5
15         # replay memory
16         self.P = deque(maxlen=self.replay_memory_size)
17
18         # model
19         self.init_model()
20
21         # variables
22         self.current_loss = 0.0
23
24
25     def init_model(self):
26         # (0) 入力文字列
27         self.x = tf.placeholder(tf.float32, [None, self.N],name='x')
28
29         # (1) サイズ変更
30         x1 = tf.reshape(self.x, [-1,self.N])
31
32         # (2) 隠れ層1
33         # 重みとバイアス
34         self.w1 = tf.Variable(tf.truncated_normal([self.N,self.M],stddev=0.01))
35         self.b1 = tf.Variable(tf.zeros([self.M]))
36         x2 = tf.nn.relu( tf.matmul(x1, self.w1) + self.b1)
37
38         # (3) 隠れ層2
39         # 重みとバイアス
40         self.w2 = tf.Variable(tf.truncated_normal([self.M,self.L],stddev=0.01))
41         self.b2 = tf.Variable(tf.zeros([self.L]))
42         x3 = tf.nn.relu( tf.matmul(x2,self.w2) + self.b2)
43
44         # (3) 出力層
45         # 重みとバイアス
46         self.w3 = tf.Variable(tf.truncated_normal([self.L,self.Nact],stddev=0.01))
47         self.b3 = tf.Variable(tf.zeros([self.Nact]))
48         self.Qs = tf.matmul(x3, self.w3) + self.b3
49
50
51         ##target network
52         # (0) 入力文字列
53         self.xtar = tf.placeholder(tf.float32, [None, self.N],name='x')
54
55         # (1) サイズ変更
56         x1 tar = tf.reshape(self.xtar, [-1,self.N])
57
58         # (2) 隠れ層1
59         # 重みとバイアス

```

```

60 self.w1tar = tf.Variable(self.w1)
61 self.b1tar = tf.Variable(self.b1)
62 x2tar = tf.nn.relu( tf.matmul(x1tar, self.w1tar) + self.b1tar)
63
64 # (3) 隠れ層2
65 # 重みとバイアス
66 self.w2tar = tf.Variable(self.w2)
67 self.b2tar = tf.Variable(self.b2)
68 x3tar = tf.nn.relu( tf.matmul(x2tar, self.w2tar) + self.b2tar)
69
70 # (3) 出力層
71 # 重みとバイアス
72 self.w3tar = tf.Variable(self.w3)
73 self.b3tar = tf.Variable(self.b3)
74 self.Qstar = tf.matmul(x3tar, self.w3tar) + self.b3tar
75
76
77 # 損失関数
78 self.y = tf.placeholder(tf.float32, [None])#バッチサイズ分の長さをもつ
79 self.act_indices=tf.placeholder(tf.int32, [None,2])#選択アクションを含むindex
80 self.Q= tf.gather_nd(self.Qs,self.act_indices)#バッチサイズ分の長さをもつQ値
81
82 self.loss = tf.reduce_mean(tf.losses.huber_loss(self.y,self.Q))
83
84 # train operation
85 optimizer = tf.train.AdamOptimizer()
86 self.training = optimizer.minimize(self.loss)
87
88 # session
89 self.sess = tf.Session()
90 self.sess.run(tf.global_variables_initializer())
91
92 def Qs_func(self, state):
93     res=self.sess.run(self.Qs, feed_dict={self.x: [state]})
94     return res[0]
95
96 def Qstar_func(self,state):
97     res=self.sess.run(self.Qstar, feed_dict={self.xtar: [state]})
98     return res[0]
99
100 def target_update(self):
101     tf.assign(self.w1tar,self.w1)
102     tf.assign(self.b1tar, self.b1)
103     tf.assign(self.w2tar, self.w2)
104     tf.assign(self.b2tar, self.b2)
105     tf.assign(self.w3tar, self.w3)
106     tf.assign(self.b3tar, self.b3)
107
108 def select_action(self, state):
109     if np.random.rand() <= self.epsilon:
110         # random
111         return np.random.randint(self.Nact)
112     else:
113         # max_action Q(state, action)
114         return np.argmax(self.Qs_func(state))
115
116 def store_experience(self, state, action, reward, state_1, terminal):
117     self.P.append((state, action, reward, state_1, terminal))
118
119 def experience_replay(self):
120     state_minibatch = []

```

```

121 y_minibatch = []
122 act_indices=[]
123 index=0
124
125 # sample random minibatch
126 # キューに十分にたまるまでは全部使う
127 minibatch_size = min(len(self.P), self.minibatch_size)
128 minibatch_indexes = np.random.choice(np.arange(len(self.P)), minibatch_size,replace=False)
129
130 for j in minibatch_indexes:
131     state_j, action_j, reward_j, state_j_1, terminal = self.P[j]
132
133     y_j= 0
134
135     if terminal:
136         y_j= reward_j
137     else:
138         y_j= reward_j + np.max(self.Qstar_func(state_j_1))
139
140     state_minibatch.append(state_j)
141     y_minibatch.append(y_j)
142     act_indices.append([index,action_j])
143     index+=1
144
145 # training
146 self.sess.run(self.training, feed_dict={self.x: state_minibatch, self.y: y_minibatch, self.act_in
147
148 # for log
149 self.current_loss = self.sess.run(self.loss, feed_dict={self.x: state_minibatch, self.y: y_miniba

```

In [148]:

```

1 class StateArray:
2     def __init__(self,N):
3         self.N=N
4         self.arr=np.random.randint(0,9,N)
5         self.frame=0
6
7     def excute_action(self,action):
8         temp=self.arr[action]
9         self.arr[action]=self.arr[action+1]
10        self.arr[action+1]=temp
11
12    def is_sorted(self):
13        res=True
14        for i in range(self.N-1):
15            res=(self.arr[i]<=self.arr[i+1]) and res
16        return res
17
18    def observe(self):
19        self.reward=-1*2/(self.N*(self.N-1))
20        self.terminal=False
21        if self.is_sorted():
22            self.reward=1
23            self.terminal=True
24        return self.arr,self.reward,self.terminal

```

In [149]:

```
1  # parameters
2  n_epochs = 100
3
4  decrease=0.001**(1/n_epochs)
5
6  N=4
7
8  dqn = DQN_SORT(N)
9
10 A=1024
11 B=1024
12
13 count=0
14
15 learning_curve=[]
16 average=0.0
17 buff=3
18
19 for e in range(n_epochs):
20     # reset
21     frame=0
22     cum_reward=0.0
23     sa=StateArray(N)
24     state_t_1,reward_t,terminal=sa.observe()
25
26     while not terminal:
27         state_t = state_t_1
28
29         # execute action in environment
30         action_t = dqn.select_action(state_t)
31         sa.excute_action(action_t)
32
33         # observe environment
34         state_t_1, reward_t, terminal = sa.observe()
35
36         # store experience
37         dqn.store_experience(state_t, action_t, reward_t, state_t_1, terminal)
38
39         # experience replay
40         if count>A:
41             dqn.experience_replay()
42
43         if count%B==0:
44             dqn.target_update()
45
46
47         cum_reward+=reward_t
48         average+=cum_reward
49
50         if frame%buff==buff-1:
51             average/=buff
52             learning_curve.append([count,average])
53
54         count+=1
55         frame+=1
56
57     dqn.epsilon*=decrease
58     if e%10==0 and count>A:
59         print("epoch=",e,"loss=",dqn.current_loss)
```

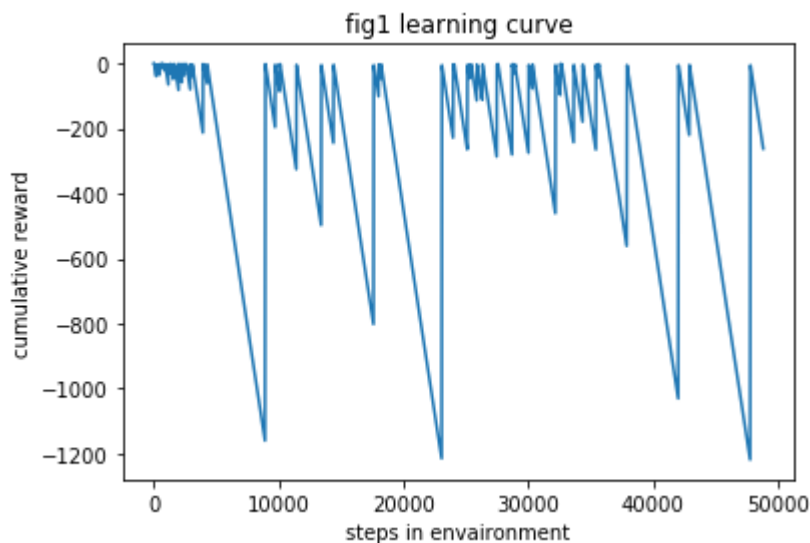
```
60  
61 print("finished")  
62 learning_curve=np.array(learning_curve)  
63 learning_curve=learning_curve.T
```

```
epoch= 20 loss= 0.00044829276  
epoch= 30 loss= 3.8358714e-05  
epoch= 40 loss= 0.0006216031  
epoch= 50 loss= 9.5349577e-07  
epoch= 60 loss= 6.2366175e-06  
epoch= 70 loss= 2.3341701e-05  
epoch= 80 loss= 3.7558806e-05  
epoch= 90 loss= 1.4010259e-05  
finished
```

Plot Learning Curve

In [150]:

```
1 plt.title("fig1 learning curve")  
2 plt.xlabel("steps in envaironment")  
3 plt.ylabel("cumulative reward")  
4 plt.plot(learning_curve[0],learning_curve[1])  
5 plt.show()
```



Conclusion

急激に値が上昇している点はエポックの移り変わりを示すが、むしろ学習初期の方がエポック間のstep数が少ない(=ソートが早く終わっている)ことがわかる。学習がうまくいっていないようである。残念ながらレポート提出までに適切なパラメータを見つけることができなかった。

Advance

N-1個の行動に加えて、自らタスクを終了させる行動を加え、終了時点で正しくソートできていれば報酬を+1、できていなければ報酬を+0。

In [117]:

```

1  class DQN_SORT_ADV:
2      def __init__(self,inputnum):
3          self.N=inputnum
4          #隠れ層1のニューロン数
5          self.M=100
6          #隠れ層2のニューロン数
7          self.L=100
8          #行動数
9          self.Nact=self.N
10         #ミニバッチ数
11         self.minibatch_size = 32
12         #キューのサイズ
13         self.replay_memory_size = 10000
14         self.epsilon=0.8
15         # replay memory
16         self.P = deque(maxlen=self.replay_memory_size)
17
18         # model
19         self.init_model()
20
21         # variables
22         self.current_loss = 0.0
23
24
25     def init_model(self):
26         # (0) 入力文字列
27         self.x = tf.placeholder(tf.float32, [None, self.N],name='x')
28
29         # (1) サイズ変更
30         x1 = tf.reshape(self.x, [-1,self.N])
31
32         # (2) 隠れ層1
33         # 重みとバイアス
34         self.w1 = tf.Variable(tf.truncated_normal([self.N,self.M],stddev=0.01))
35         self.b1 = tf.Variable(tf.zeros([self.M]))
36         x2 = tf.nn.relu( tf.matmul(x1, self.w1) + self.b1)
37
38         # (3) 隠れ層2
39         # 重みとバイアス
40         self.w2 = tf.Variable(tf.truncated_normal([self.M,self.L],stddev=0.01))
41         self.b2 = tf.Variable(tf.zeros([self.L]))
42         x3 = tf.nn.relu( tf.matmul(x2,self.w2) + self.b2)
43
44         # (3) 出力層
45         # 重みとバイアス
46         self.w3 = tf.Variable(tf.truncated_normal([self.L,self.Nact],stddev=0.01))
47         self.b3 = tf.Variable(tf.zeros([self.Nact]))
48         self.Qs = tf.matmul(x3, self.w3) + self.b3
49
50
51         ##target network
52         # (0) 入力文字列
53         self.xtar = tf.placeholder(tf.float32, [None, self.N],name='x')
54
55         # (1) サイズ変更
56         x1 tar = tf.reshape(self.xtar, [-1,self.N])
57
58         # (2) 隠れ層1
59         # 重みとバイアス

```



```

60 self.w1tar = tf.Variable(self.w1)
61 self.b1tar = tf.Variable(self.b1)
62 x2tar = tf.nn.relu( tf.matmul(x1tar, self.w1tar) + self.b1tar)
63
64 # (3) 隠れ層2
65 # 重みとバイアス
66 self.w2tar = tf.Variable(self.w2)
67 self.b2tar = tf.Variable(self.b2)
68 x3tar = tf.nn.relu( tf.matmul(x2tar, self.w2tar) + self.b2tar)
69
70 # (3) 出力層
71 # 重みとバイアス
72 self.w3tar = tf.Variable(self.w3)
73 self.b3tar = tf.Variable(self.b3)
74 self.Qstar = tf.matmul(x3tar, self.w3tar) + self.b3tar
75
76
77 # 損失関数
78 self.y = tf.placeholder(tf.float32, [None])#バッチサイズ分の長さをもつ
79 self.act_indices=tf.placeholder(tf.int32, [None,2])#選択アクションを含むindex
80 self.Q= tf.gather_nd(self.Qs,self.act_indices)#バッチサイズ分の長さをもつQ値
81
82 self.loss = tf.reduce_mean(tf.losses.huber_loss(self.y,self.Q))
83
84 # train operation
85 optimizer = tf.train.AdamOptimizer()
86 self.training = optimizer.minimize(self.loss)
87
88 # session
89 self.sess = tf.Session()
90 self.sess.run(tf.global_variables_initializer())
91
92 def Qs_func(self, state):
93     res=self.sess.run(self.Qs, feed_dict={self.x: [state]})
94     return res[0]
95
96 def Qstar_func(self,state):
97     res=self.sess.run(self.Qstar, feed_dict={self.xtar: [state]})
98     return res[0]
99
100 def target_update(self):
101     tf.assign(self.w1tar,self.w1)
102     tf.assign(self.b1tar, self.b1)
103     tf.assign(self.w2tar, self.w2)
104     tf.assign(self.b2tar, self.b2)
105     tf.assign(self.w3tar, self.w3)
106     tf.assign(self.b3tar, self.b3)
107
108 def select_action(self, state):
109     if np.random.rand() <= self.epsilon:
110         # random
111         return np.random.randint(self.Nact)
112     else:
113         # max_action Q(state, action)
114         return np.argmax(self.Qs_func(state))
115
116 def store_experience(self, state, action, reward, state_1, terminal):
117     self.P.append((state, action, reward, state_1, terminal))
118
119 def experience_replay(self):
120     state_minibatch = []

```

```
121 y_minibatch = []
122 act_indices=[]
123 index=0
124
125 # sample random minibatch
126 # キューに十分にたまるまでは全部使う
127 minibatch_size = min(len(self.P), self.minibatch_size)
128 minibatch_indexes = np.random.choice(np.arange(len(self.P)), minibatch_size,replace=False)
129
130 for j in minibatch_indexes:
131     state_j, action_j, reward_j, state_j_1, terminal = self.P[j]
132
133     y_j= 0
134
135     if terminal:
136         y_j= reward_j
137     else:
138         y_j= reward_j + np.max(self.Qstar_func(state_j_1))
139
140     state_minibatch.append(state_j)
141     y_minibatch.append(y_j)
142     act_indices.append([index,action_j])
143     index+=1
144
145 # training
146 self.sess.run(self.training, feed_dict={self.x: state_minibatch, self.y: y_minibatch, self.act_in
147
148 # for log
149 self.current_loss = self.sess.run(self.loss, feed_dict={self.x: state_minibatch, self.y: y_miniba
```

In [118]:

```
1 class StateArrayAdv:
2     def __init__(self,N):
3         self.N=N
4         self.arr=np.random.randint(0,9,N)
5         self.frame=0
6         self.reward=-1*2/(self.N*(self.N-1))
7         self.terminal=False
8
9     def excute_action(self,action):
10        if action!=self.N-1:
11            temp=self.arr[action]
12            self.arr[action]=self.arr[action+1]
13            self.arr[action+1]=temp
14        else:
15            self.terminal=True
16
17    def is_sorted(self):
18        res=True
19        for i in range(self.N-1):
20            res=(self.arr[i]<=self.arr[i+1]) and res
21        return res
22
23    def observe(self):
24        self.reward=-1*2/(self.N*(self.N-1))
25        if self.terminal:
26            if self.is_sorted():
27                self.reward=1
28            else:
29                self.reward=0
30
31        return self.arr,self.reward,self.terminal
```

In [119]:

```

1  n_epochs=5000
2  dqn_adv = DQN_SORT_ADV(N)
3
4  count=0
5
6  learning_curve_adv=[]
7  average=0.0
8
9  N=4
10
11 A=1024
12 B=1024
13
14 for e in range(n_epochs):
15     # reset
16     frame=0
17     cum_reward=0.0
18     sa=StateArrayAdv(N)
19     state_t_1,reward_t,terminal=sa.observe()
20
21     while not terminal:
22         state_t = state_t_1
23
24         # execute action in environment
25         action_t = dqn_adv.select_action(state_t)
26         sa.excute_action(action_t)
27
28         # observe environment
29         state_t_1, reward_t, terminal = sa.observe()
30
31         # store experience
32         dqn_adv.store_experience(state_t, action_t, reward_t, state_t_1, terminal)
33
34         # experience replay
35         if count>A:
36             dqn_adv.experience_replay()
37
38
39         if count%B==0:
40             dqn_adv.target_update()
41
42
43         cum_reward+=reward_t
44         average+=cum_reward
45
46         if frame%buff==buff-1:
47             average/=buff
48             learning_curve_adv.append([count,average])
49
50         count+=1
51         frame+=1
52
53     dqn_adv.epsilon*=decrease
54     if e%100==0 and count>A:
55         print("epoch=",e,"loss=",dqn_adv.current_loss)
56
57     print("finished")
58     learning_curve_adv=np.array(learning_curve_adv)
59     learning_curve_adv=learning_curve_adv.T

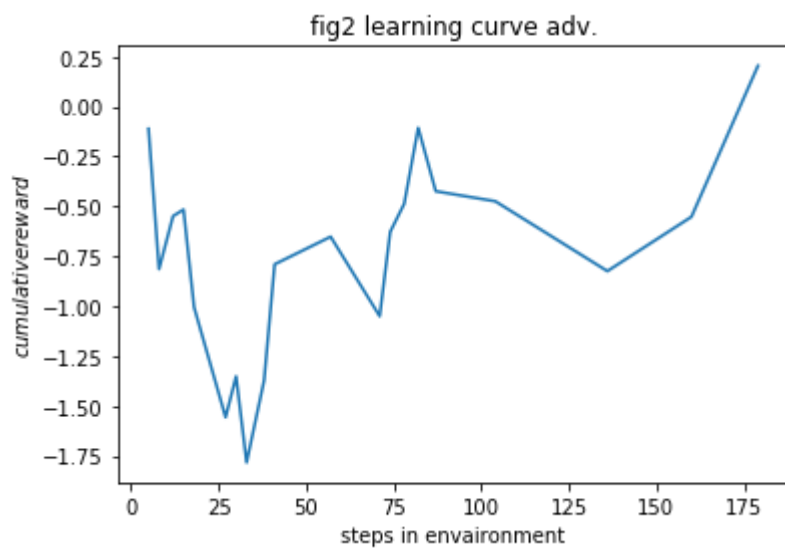
```

```
epoch= 1000 loss= 0.0015185372
epoch= 1100 loss= 0.00039484084
epoch= 1200 loss= 0.01871253
epoch= 1300 loss= 0.012573384
epoch= 1400 loss= 0.00016127527
epoch= 1500 loss= 0.021034678
epoch= 1600 loss= 0.026605979
epoch= 1700 loss= 0.015753468
epoch= 1800 loss= 0.014635619
epoch= 1900 loss= 7.38577e-06
epoch= 2000 loss= 0.0016783845
epoch= 2100 loss= 0.009081308
epoch= 2200 loss= 6.0189937e-05
epoch= 2300 loss= 0.014398257
epoch= 2400 loss= 0.0062923785
epoch= 2500 loss= 0.00012045413
epoch= 2600 loss= 0.0075981985
epoch= 2700 loss= 3.6846986e-06
epoch= 2800 loss= 0.00081098726
epoch= 2900 loss= 9.89033e-06
epoch= 3000 loss= 0.0011454413
epoch= 3100 loss= 0.00040335642
epoch= 3200 loss= 0.0006826106
epoch= 3300 loss= 0.002611473
epoch= 3400 loss= 0.00024435957
epoch= 3500 loss= 0.00018714493
epoch= 3600 loss= 0.0001550679
epoch= 3700 loss= 9.0620296e-07
epoch= 3800 loss= 4.152867e-05
epoch= 3900 loss= 7.973924e-05
epoch= 4000 loss= 1.8670858e-05
epoch= 4100 loss= 1.2622099e-06
epoch= 4200 loss= 0.00012013019
epoch= 4300 loss= 0.00052570656
epoch= 4400 loss= 9.479394e-06
epoch= 4500 loss= 1.4638972e-06
epoch= 4600 loss= 5.8158326e-05
epoch= 4700 loss= 2.265943e-05
epoch= 4800 loss= 2.5988376e-05
epoch= 4900 loss= 2.4094043e-06
finished
```

Plot Learning Curve

In [120]:

```
1 plt.title("fig2 learning curve adv.")
2 plt.xlabel("steps in envaironment")
3 plt.ylabel("$cumulative reward$")
4 plt.plot(learning_curve_adv[0],learning_curve_adv[1])
5 plt.show()
```



Conclusion

今度は学習がうまくいってる。また、終了することでstep数が抑えられ、非常に早く学習が終了する。step数を抑えることで、累積したマイナスの報酬が抑えられ、成功報酬の+1が有効に機能している。終了行動のない場合は累積したマイナスの報酬があまりにも大きな値になり、成功報酬+1と比較して相対的にほとんど効果をなしていなかった。その点こちらはこの効果が反映され学習にいい影響を及ぼしていると考えられる。