

TP2, Reinforcement learning

Mariia Drozdova

February 25, 2020

1 Exercise 0

File `replacement_model.py` contains a model of deterioration.

1.1 0.1 Identify a law of deterioration between two decisions, deterioration when the object should be replaced and associated punishment. Interpret a function of maintenance cost.

```
1 xmax = 10 # maximum level of use
  xdead = 9 # level of death
3 Creplace = 40 # cost of replacement
  Cdead = 10; # cost of death
```

Listing 1: Parameters in `replacement_model.py`

First, we notice that **x** here states for current level of use. When it is equal to **xdead**, the object should be recycled that has a cost **Cdead** (punishment). Buying new object will instead will cost **Creplace**. If the object is replaced when it is not dead, we do not have to pay **Cdead**.

```
beta = 2.0 # exponential distribution parameter
2
def sample_exp( beta ):
4     return -1.0/beta*log( random() )

6 def next_state_and_reward( x, a ):
    c=0
8     if a==0: # replace
        y=sample_exp(beta)
10        c+=Creplace + maintenance_cost(0,y)
    else: # keep
12        y = x + sample_exp(beta)
        c += maintenance_cost(x, y)
14    if y>xdead:
        y, c2 = next_state_and_reward ( y, 0 )
16        c += Cdead + c2
    return y, c
```

Listing 2: The law for the level of use

The level of use which corresponds to the level of deterioration is a sum of values produced by `sample_exp`. This function returns samples from the exponential distribution of parameter beta. Sum of exponential variables with the same beta is gamma. A law of deterioration between two decisions is exponential.

```

1 def maintenance_cost(x,y):
    c=y-x
3     if x < 2.5 < y:
        c=c+10
5     if x < 5 < y:
        c=c+20
7     if x < 7.5 < y:
        c=c+30
9     return c

```

Listing 3: Maintenance cost function.

The maintenance cost can be interpreted as the following:

The cost is equal to the change of the level of use. If the level of use passes some predefined points, addition service should be applied. For the point of 2.5, the service costs 10, for 5 the cost is 20 and, finally, for 7.5 the cost is 30.

1.2 How to improve the solution?

A file `replacement.py` contains an already implemented algorithm of Discretized Value Iteration. We discretize the value space - the possible values of the level of use **x**. Then, we compute all possible transitions and rewards from this **x** depending on the taken action (**y0** and **y1**).

The file contains the following parameters:

```

1 N = 151 # number of samples points
M = 10 # number of next_state_and_reward samples
3 K = 30 # number of updating alpha

```

Listing 4: Maintenance cost function.

N is responsible for discretisizing.

We construct the value function which associates each sample interval to the value which can be received from it. The algorithm updates these values though **K** iterations. For each iteration we run the algorithm for all sample points **N** making **M** next steps from each.

The figures 1 and 2 are the example of the output we get if we run only preimplemented algorithm.

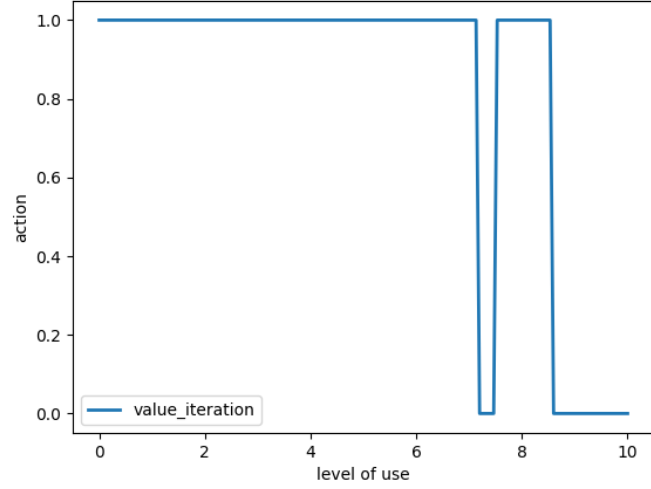
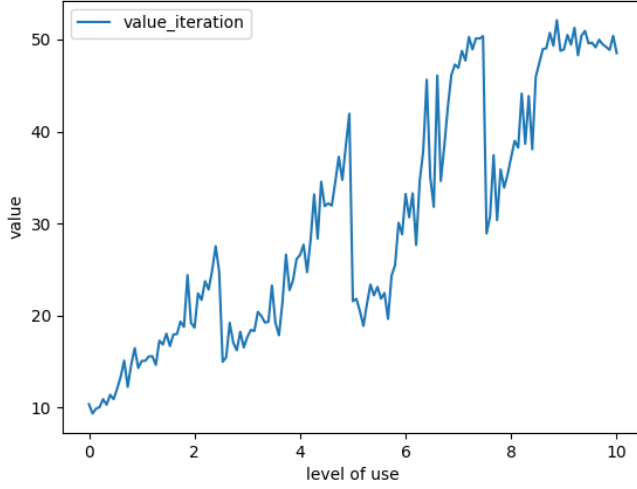


Figure 1: Values for the Discretized Value Iteration. Figure 2: Actions for the Discretized Value Iteration.

Thus, to improve the solution we can use more discretized space. (bigger \mathbf{N}) We can also augment the number of next steps (\mathbf{M}) performed for each sample. The higher number of step leads to more precision.

2 Exercice 1 (Fitted Value-Iteration)

In this exercise we try to approximate the value function in the Fourier basis of functions $\phi_j(x) = \cos \frac{\pi j x}{x_{\max}}$ where $1 \leq j \leq d$. This time we will update the coefficients of the value function in this basis through solving minimization problem.

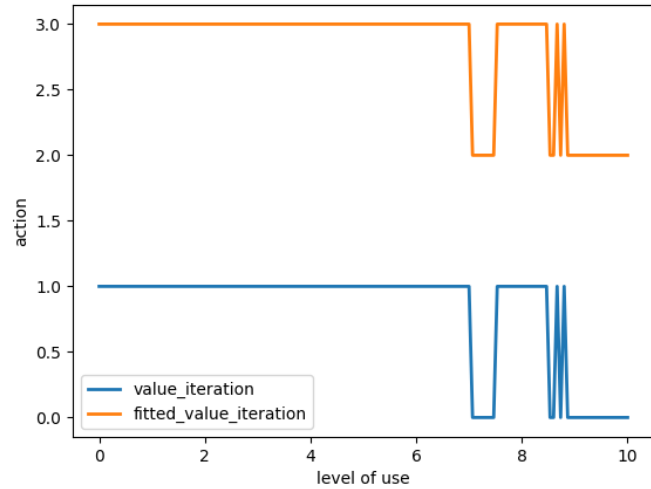
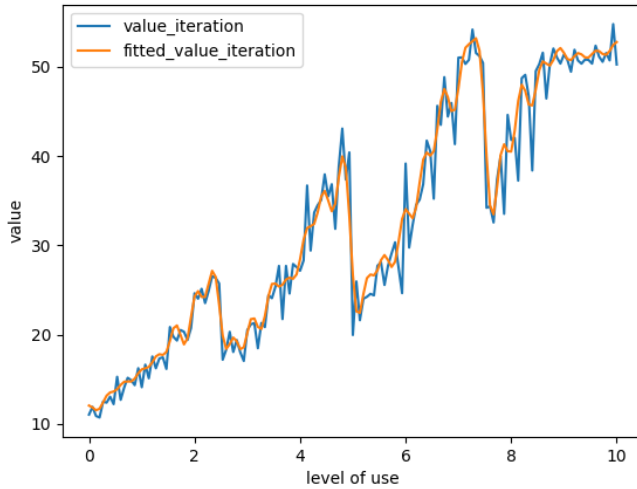


Figure 3: Values for the Fitted Value Iteration.

Figure 4: Taken actions for the Fitted Value Iteration.

The Figures 3 and 4 show the approximated value function and the action of this and previous technique. The value function is smooth in case of fitted value iteration as we have smooth basis. In terms of basis the first algorithm of value iteration can be rewritten as if we have a discrete euclidean basis of N vectors in N -dimensional space. The coefficients are corresponding values of each of N samples.

```

1 def return_cos(j, x, x_max=xmax):
2     return np.cos(j*np.pi*x/x_max)
3 def compute_F(x):
4     F = []
5     for j in range(d):
6         F.append(return_cos(j, x))
7     F = np.array(F)
8     return F
9 def return_value(alpha, x):
10    res = 0
11    for j in range(len(alpha)):
12        res += alpha[j] * return_cos(j, x)
13    return res
14 Fi = compute_F(x_samples).T
15
16 def fitted_value_iterate(alpha):
17     alpha2, pol = [], []
18     Tvs = []
19     for i in range(N):
20         print((i+1), "/", N, " \r", end=' ')
21         sys.stdout.flush()
22         Tv0, Tv1 = 0, 0
23         for j in range(M):
24             value0 = return_value(alpha, y0_samples[i][j])
25             value1 = return_value(alpha, y1_samples[i][j])
26             Tv0 += r0_samples[i][j] + gamma * value0
27             Tv1 += r1_samples[i][j] + gamma * value1
28         if Tv0 <= Tv1:
29             Tv = Tv0 / M
30             pol.append(0)
31         else:
32             Tv = Tv1 / M
33             pol.append(1)
34         Tvs.append(Tv)
35     Tvs = np.array(Tvs)
36     alpha2 = np.linalg.inv((np.transpose(Fi) @ Fi)) @ (np.transpose(Fi) @ Tvs)
37     return alpha2, pol
38
39 def fitted_value_iteration(K):
40     alpha = initiate_alpha()
41     for i in range(K):
42         alpha, pol = fitted_value_iterate(alpha)
43         print(i, " ", pol)
44     values = return_value(alpha, x_samples)
45     return values, pol

```

Listing 5: Fitted Value Iteration algorithm.

3 Exercise 2 (Fitted Q-Iteration)

In this exercise we approximate the \mathbf{q} in the basis ϕ introduced in the previous section. We introduce two set of coefficients α_0 and α_1 corresponding to the two actions. Then, we solve two problems of minimization.

After getting q_0 and q_1 for each sample we find the minimal of them getting a value and an action for each sample.

The code can be found in the Listing 6. The output for all algorithms included into this project is displayed in the Figures 5 and 6. The curves for value are approximating the same law.

```

1 def fitted_q_iterate(alpha0, alpha1):
2     pol=[]
3     Tvs0, Tvs1 = [],[]
4     for i in range(N):
5         #print((i+1),"/",N,"    \r", end=' ')
6         sys.stdout.flush()
7         Tv0,Tv1=0,0
8         for j in range(M):
9             min_q0 = np.min(np.array([return_value(alpha0, y0_samples[i][j]),
10             return_value(alpha1, y0_samples[i][j])]))
11             min_q1 = np.min(np.array([return_value(alpha0, y1_samples[i][j]),
12             return_value(alpha1, y1_samples[i][j])]))
13             Tv0 += r0_samples[i][j] + gamma * min_q0
14             Tv1 += r1_samples[i][j] + gamma * min_q1
15             Tvs0.append(Tv0/M)
16             Tvs1.append(Tv1/M)
17         Tvs0 = np.array(Tvs0).reshape(N, 1)
18         Tvs1 = np.array(Tvs1).reshape(N, 1)
19         alpha0_2 = np.linalg.inv((np.transpose(Fi) @ Fi) @ (np.transpose(Fi) @
20         Tvs0))
21         alpha1_2 = np.linalg.inv((np.transpose(Fi) @ Fi) @ (np.transpose(Fi) @
22         Tvs1))
23         return alpha0_2, alpha1_2, pol
24
25 def fitted_q_iteration(K):
26     alpha0, alpha1=initiate_alpha(), initiate_alpha()
27     for i in range(K):
28         alpha0, alpha1, pol = fitted_q_iterate(alpha0, alpha1)
29         print(i,"    ")# ,pol
30
31     qa0 = [return_value(alpha0, x_sample) for x_sample in x_samples]
32     qa1 = [return_value(alpha1, x_sample) for x_sample in x_samples]
33     values = np.min(np.stack((qa0, qa1)), axis=0)
34     pol = np.argmin(np.stack((qa0, qa1)), axis=0)
35     return values, pol

```

Listing 6: Fitted Q Iteration algorithm.

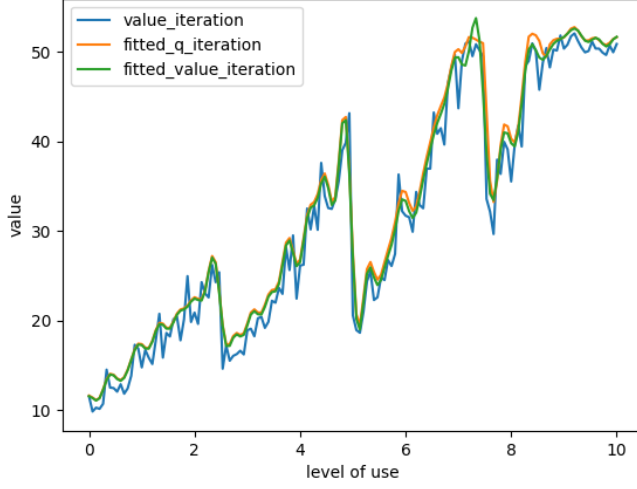


Figure 5: Values depending on the level of use.

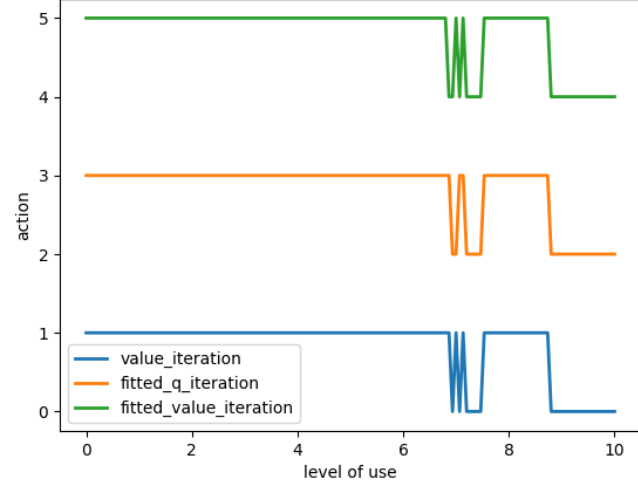


Figure 6: Actions depending on the level of use.

Finally we test the three algorithms with M equal to one which means that we generate next step only once. It leads to big uncertainties for all algorithms with the worse result for discretized version. The difference between fitted value and q iteration is also becoming more visible. In case of value iteration we approximate two sets of coefficients which leads to bigger variance and error. Most of the time the value returned for q iteration is higher than the one returned for value iteration. Q-iteration has worse convergence properties than value iteration in our case. If we take bigger M , on the contrary the curves are getting closer to each other.

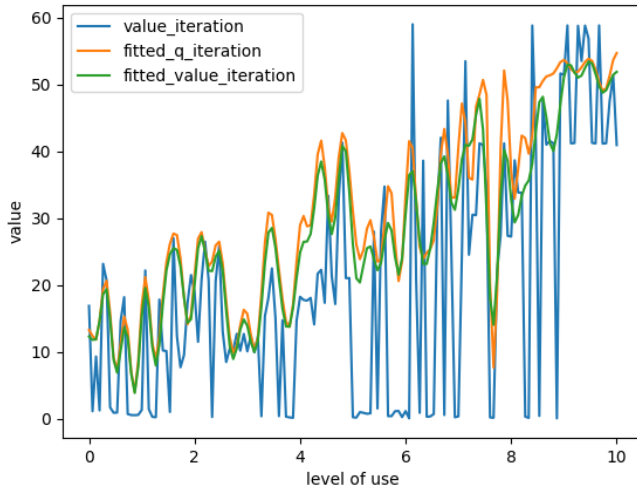


Figure 7: Values depending on the level of use with $M=1$.

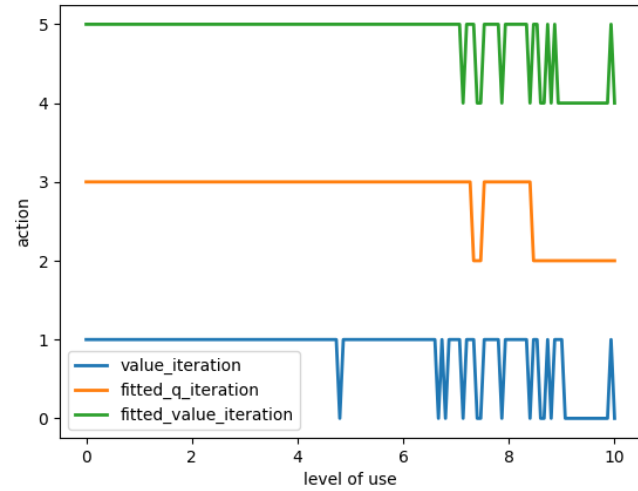


Figure 8: Actions depending on the level of use $M=1$.