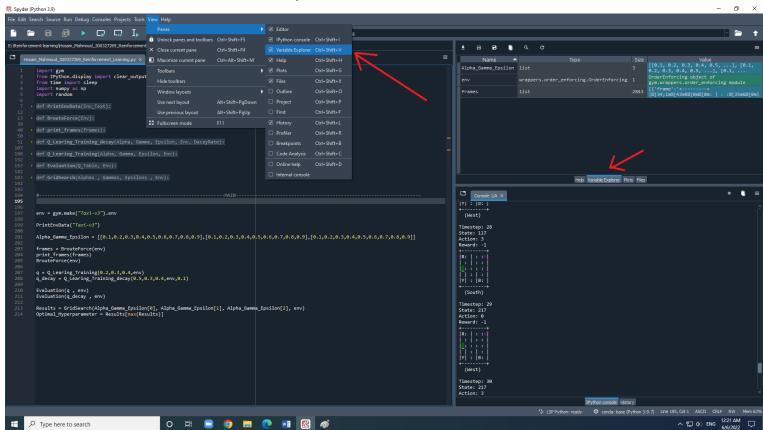For best experience please consider to run **Hosam_Mahmoud_300327269_Reinforcement_Learning.py** on Spyder to see all created variables on variables explorer tab.
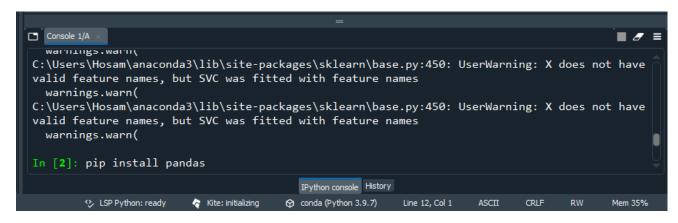
Please make sure that you have installed these libraries before importing it ->
*import gym*
*from IPython.display import clear_output*
*from time import sleep*
*import numpy as np*
*import random*

------------------------------------------------------------------------

To install the missing libraries, you have to write something like this ->

>>> **pip install** *numpy*

In Spyder console.



or

>>> **conda install pandas**

In Anaconda prompt.

Here I have 7 functions in my code.

**def PrintEnvData(Env_Text):** is for printing the number of actions and states of the environments.

**def BrouteForce(Env):** to try all the possible paths for one state.

**def print_frames(frames):** to animate the paths.

**def Q_Learing_Training_decay(Alpha, Gamma, Epsilon, Env, DecayRate):** Training function using a decay over episodes.

Here is the parameter tuning part, but in case of **gamma** and **epsilon** I have tried to change its values but I got an infinite loop when applying evaluation.

```python
            if i % 100 == 0:
                clear_output(wait = True)
                print(f"Episode: {i}")
                alpha = abs(alpha - (1/(1 + (DecayRate * 100000))) * alpha)
                # gamma = abs(gamma - (1/(1 + (DecayRate * 100000))) * gamma)
                # epsilon = abs(epsilon - (1/(1 + (DecayRate * 100000))) * epsilon)

                print('Alpha = ', alpha)
                print('Gamma = ', gamma)
                print('Epsilon = ', epsilon)

                alpha = Alpha if alpha == 0 else alpha
                # gamma = Gamma if gamma == 0 else gamma
                # epsilon = Epsilon if epsilon == 0 else epsilon

        print("Training finished.\n")
        return q_table
```

**def Q_Learing_Training(Alpha, Gamma, Epsilon, Env):** Training function without parameter tuning.

**def Evaluation(Q_Table, Env):** to return the performance indicator of each reinforcement model.

**def GridSearch(Alphas , Gammas, Epsilons , Env):** to apply GridSearch and optain the optimal hyperparameters.
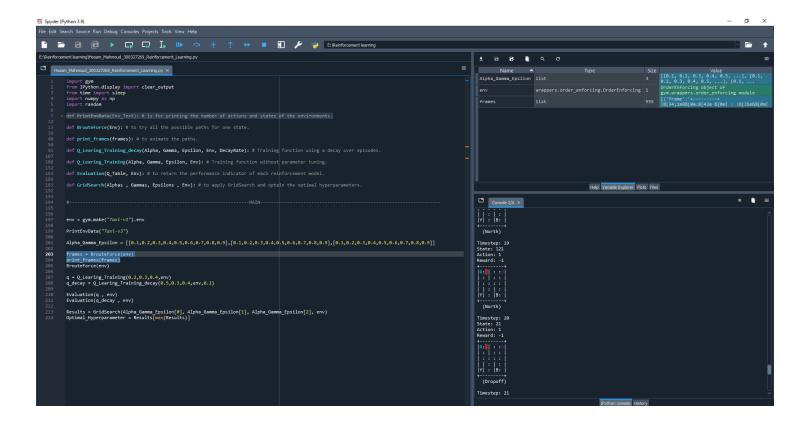
```
193
194    #-------------------------------------------------MAIN-------------------------------------------------
195
196
197    env = gym.make("Taxi-v3").env
198
199    PrintEnvData("Taxi-v3")
200
201    Alpha_Gamma_Epsilon = [[0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9],[0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9],[0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9]]
202
203    frames = BrouteForce(env)
204    print_frames(frames)
205    BrouteForce(env)
206
207    q = Q_Learing_Training(0.2,0.3,0.4,env)
208    q_decay = Q_Learing_Training_decay(0.5,0.3,0.4,env,0.1)
209
210    Evaluation(q , env)
211    Evaluation(q_decay , env)
212
213    Results = GridSearch(Alpha_Gamma_Epsilon[0], Alpha_Gamma_Epsilon[1], Alpha_Gamma_Epsilon[2], env)
214    Optimal_Hyperparameter = Results[max(Results)]
```

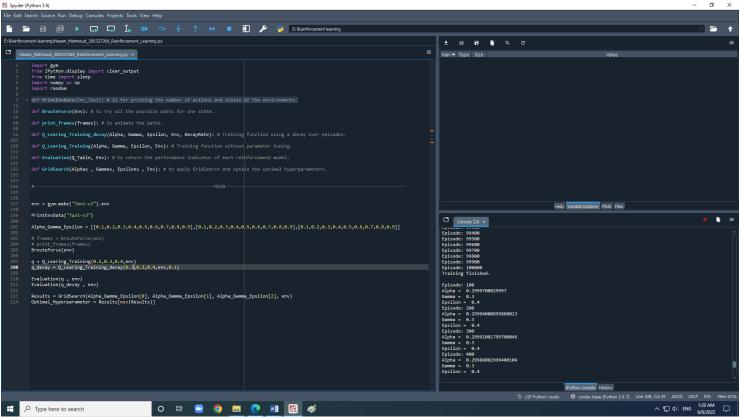For the MAIN part I have created environment and print its number of states and actions.

```
In [7]: runfile('E:/Reinforcement Learning/Hosam_Mahmoud_300327269_Reinforcement_Learning.py', wdir='E:/
Reinforcement Learning')
Action Space Discrete(6)
State Space Discrete(500)

In [8]:
```

After that I have created list of each parameter (alpha, gamma and epsilon), and I'll use it for grid search **-> 9*9*9** but this toke too long time, so I tried with another lists **-> [0.3,0.6,0.9]** which will give us a total number of combinations = 3*3*3

And after that I animated the paths of Bruteforce function.

After that I trained the two models (using decay and without using decay).
And the values of alpha are changed every 100 episodes.

After Evaluation (without using decay)

```
Results after 1000 episodes:
Average timesteps per episode: 13.284
Average penalties per episode: 0.0
Average reward per episode: 20.0
```

After Evaluation (using decay)

```
Results after 1000 episodes:
Average timesteps per episode: 12.976
Average penalties per episode: 0.0
Average reward per episode: 20.0
```

And we found that (with using decay) is faster than the normal model.

After that I applied GridSearch on function (without using Decay)

```
182
183    def GridSearch(Alphas , Gammas, Epsilons , Env): # to apply GridSearch and optain the optimal hyperparameters.
184        Dictionary = {}
185        for Alpha in Alphas:
186            for Gamma in Gammas:
187                for Epsilon in Epsilons:
188                    q = Q_Learing_Training(Alpha, Gamma , Epsilon, Env)
189                    Metric = Evaluation(q , Env)
190                    Dictionary[Metric] = [Alpha, Gamma , Epsilon]
191        return Dictionary
192
```
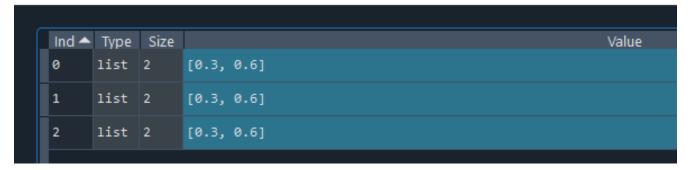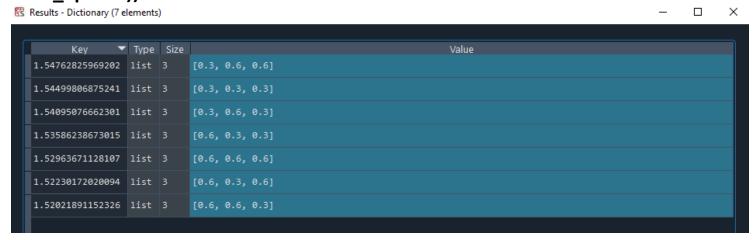
For experiments purposes I have choose small values for parameters so the program run faster -> 2*2*2

Alpha_Gamma_Epsilon - List (3 elements)

| Ind | Type | Size | Value |
|---|---|---|---|
| 0 | list | 2 | [0.3, 0.6] |
| 1 | list | 2 | [0.3, 0.6] |
| 2 | list | 2 | [0.3, 0.6] |

After applying GridSearch we will choose the maximum key and return its parameter
And the key is equal to that equation -> **Metric = (total_reward/(total_penalties + total_epochs)).**

Results - Dictionary (7 elements) — □ ✕

| Key | Type | Size | Value |
|---|---|---|---|
| 1.54762825969202 | list | 3 | [0.3, 0.6, 0.6] |
| 1.54499806875241 | list | 3 | [0.3, 0.3, 0.3] |
| 1.54095076662301 | list | 3 | [0.3, 0.6, 0.3] |
| 1.53586238673015 | list | 3 | [0.6, 0.3, 0.3] |
| 1.52963671128107 | list | 3 | [0.6, 0.6, 0.6] |
| 1.52230172020094 | list | 3 | [0.6, 0.3, 0.6] |
| 1.52021891152326 | list | 3 | [0.6, 0.6, 0.3] |

And Finally, we obtained the optimal hyperparameters.

```
Results = GridSearch(Alpha_Gamma_Epsilon[0], Alpha_Gamma_Epsilon[1], Alpha_Gamma_Epsilon[2], env)
Optimal_Hyperparameter = Results[max(Results)]
```

Optimal_Hyperparameter - List (3 elements)

| Ind | Type | Size | Value |
|---|---|---|---|
| 0 | float | 1 | 0.3 |
| 1 | float | 1 | 0.6 |
| 2 | float | 1 | 0.6 |