

Introduction to Complex problem solving

Assignment 3

INFO8003-1
Damien Ernst

Olivier Moitroux
s152122

Liège

Février 2019

Domain representation

The **origin** of the grid has been put in the **bottom left** of the board so that, in the provided domain instance of the assignment, the red cell is considered to be at position (0,1). The value of y increases as we go up (0,1) while the value of x increases by going to the right (1,0). The value of the reward `grid[0, 0]` is thus considered to be -20.

-3	1	-5	0	19
6	3	8	9	10
5	-8	4	1	-8
6	-9	4	19	-5
-20	-17	-4	-3	9

If more information relative to the implementation is expected, please refer to more details explanation in the assignment 1.

Part I

$$\gamma = 0.99$$

6 Q-learning in a Batch Setting

6.1 Convergence speed of \hat{Q} towards Q

6.1.1 Chosen framework

Assignment was not clear about the framework we had to use to display convergence. I chose this one:

- No experience replay for q-learning
- 1 single trajectory of length increasing up to convergence for q-learning.
- As many updates in q-learning as there are transitions in the generated trajectory.
- Q is computed like in question 4 with $N=500$. \hat{Q} computed by q-learning described above.

\hat{Q} is computed with an increasing trajectory length. The Q considered is always Q_{500} , that is, the Q computed in question 4 at convergence. The convergence is highlighted by considering the sum of square error between all elements of $Q(x, u)$ (that is a $width \times height \times 4$ matrix computed in Q4) and of \hat{Q} , computed with q-learning (same dimensions as $Q(x, u)$ of course). To do so, I store a $width \times height \times 4$ matrix of square error and I only update the error on the element that is updated in the q-learning. Then, I only recompute the sum for this given length of trajectory for the plot. Let's note we could also use the bound on the error like explained in the course to stop when appropriate accuracy is reached. Other metrics could be used as well.

6.1.2 Plots and discussion

The evolution of the sum of the square error between \hat{Q} and Q is reported in figure 1 for deterministic and stochastic settings.

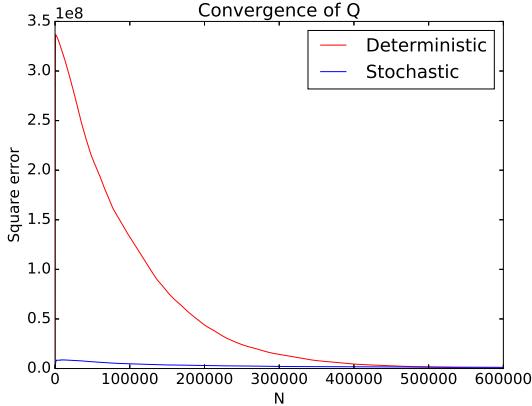


Figure 1 – Evolution of sum of square error between $Q(x, u)$ and $\hat{Q}(x, u)$, $\forall x \in X, \forall u \in U$.

At a first sight, the results are a bit surprising. Indeed, it seems like using the q-learning algorithm with a very small trajectory length leads to a huge sum of square error between \hat{Q} and Q in the deterministic case. This might be explained by the fact that Q is initialised to 0 everywhere and that the information inside the history might not even update Q at the beginning if new value are below 0 (see max in the algorithm). In the stochastic setting, the agent tends to get back to $[0, 0]$, which for me, gets a reward of -20. For this reason, the expected cumulative result is smaller in the stochastic case and thus, closer to the early stages of the q-learning algorithm. Still, I'm surprised to reach such high values.

Other than that, the q-learning algorithm seems to stabilise after 500000 updates in the deterministic case and earlier at around 400000 for the stochastic case. The convergence in the deterministic case is way steeper due to the fact that the information gain

Idea of interesting plots

These plots were not asked but could confirm our intuitions and reasoning for the next sections.

- Plotting the sum of square error between $J_{\hat{u}^*}^n$, $n = 0, \dots, N$ and $J_{u^*}^N$ of question 4.
- Highlight, like I did in previous assignment, the effect of considering several trajectories instead of a long one of increasing length.
- Make the same plots with random replay and playing with the batch-size.

6.1.3 Expected cumulative reward computed with q-learning.

Once convergence is achieved, we can derive \hat{u}^* from \hat{Q} and display $J_{\hat{u}^*}^N$. I considered $N=600000$ because q-learning is very fast anyway and 1 trajectory with no experience replay (nothing stated in the assignment about the configuration to use). Table 1 shows the optimal policy and figure 2 $J_{\hat{u}^*}^N$ besides $J_{\hat{u}^*}^N$, that was computed in question 4 of previous assignment. Let's highlight that figure 2c and 2d don't have very close values, leading us to the conclusion that the stochastic setting might be difficult to tackle with q-learning. A correct tuning of the parameters is thus important. Let's note that using the start position [3, 3] better helps in

the stochastic setting to accumulate information about states around [5, 5] (in opposition to start at [0, 0] f.e.). In the same way, considering more smaller trajectories also benefit (see my plots/justifications that have been discussed in the last assignment). This won't be shown, again, in the report but can be easily verified by tweaking values in the code.

\downarrow	\rightarrow	\rightarrow	\rightarrow	\rightarrow
\rightarrow	\rightarrow	\rightarrow	\rightarrow	\uparrow
\uparrow	\rightarrow	\uparrow	\uparrow	\uparrow
\uparrow	\rightarrow	\rightarrow	\uparrow	\uparrow
\uparrow	\rightarrow	\uparrow	\uparrow	\leftarrow

Deterministic

\downarrow	\downarrow	\downarrow	\rightarrow	\uparrow
\rightarrow	\rightarrow	\rightarrow	\rightarrow	\rightarrow
\uparrow	\uparrow	\downarrow	\downarrow	\uparrow
\uparrow	\rightarrow	\leftarrow	\leftarrow	\downarrow
\uparrow	\rightarrow	\uparrow	\uparrow	\downarrow

Stochastic, $\beta = 0.2$

Table 1 – Optimal policy $\hat{\mu}^*$ by q-learning (N=600000)

<pre>Expected cumulative reward of optimal policy: [[1767.7243 1783.6172 1807.5653 1826.4805 1826.6141] [1780.2010 1796.1281 1807.3812 1817.4734 1826.5073] [1767.7762 1781.5558 1796.2676 1807.1087 1816.9887] [1754.0511 1774.3254 1788.8675 1788.8108 1789.7836] [1741.4957 1751.8190 1774.2978 1788.8282 1766.8551]]</pre>	<pre>[[1831.52989208 1846.58327896 1870.39327896 1889.39327896 1889.39327896] [1843.96958796 1859.67237896 1870.48327896 1880.39327896 1889.39327896] [1831.42382487 1844.96958796 1859.67237896 1870.48327896 1880.39327896] [1818.00351941 1838.40312487 1853.03958796 1852.67237896 1853.48327896] [1805.71741701 1815.91302641 1838.40312487 1853.03958796 1831.50919208]]</pre>
(a) Deterministic: $J_{\hat{\mu}^*}^N$	(b) Deterministic: $J_{\mu^*}^N$
<pre>Expected cumulative reward of optimal policy: [[188.2131 189.5997 197.5258 197.7263 194.3334] [192.3373 198.5971 199.0298 196.7170 196.1319] [188.2863 194.5730 198.4449 212.6047 196.8181] [189.6406 199.5346 202.3155 201.3929 202.9487] [191.7618 186.3829 199.8955 210.1213 200.3684]]</pre>	<pre>[[294.6329785 301.13313707 314.29031983 329.49031983 329.49031983] [301.12669759 309.83801903 315.78791983 322.29031983 329.49031983] [296.62553902 301.92569759 309.83801903 315.76791983 322.29031983] [292.26141443 301.65046062 309.49049759 304.23801903 309.49793707] [289.60316165 289.04115133 301.65046062 309.49049759 296.05790011]]</pre>
(c) Stochastic: $J_{\hat{\mu}^*}^N$ ($\beta=0.2$)	(d) Stochastic: $J_{\mu^*}^N$ ($\beta = 0.2$)

Figure 2 – Left: $J_{\hat{\mu}^*}^N$, computed with q-learning ($\alpha = 0.05$, no experience replay, N=600000, 1 trajectory of length N.). Right: $J_{\mu^*}^N$ (N=500), see prev. assignment Q4.

6.2 Intelligent agent with ϵ -greedy policy.

6.2.1 Analysis of effect of batch size

Framework:

- Constants:
 - 100 episodes
 - 1000 trajectories
 - $\iff 100 \cdot 1000 = 100000$ total transitions considered for learning.
 - $\alpha = 0.05$
- Parameter: batch-size
 - Experience replay, batchsize = 1
 - Experience replay, batchsize = 10
 - Experience replay, batchsize = 100

Results and discussion.

Changing the batch-size means we can sample i.i.d. several tuples in the form $[x_i, u_i, r_i, x_i +$

$1] i \in [1, T]$ at once. Of course, we fix the total number of sampling for the q-learning so that the only parameter is the batch-size. **By increasing the batch-size, there is a gain in information.** Indeed, by fixing the batch-size to x , we update directly x times Q and we are guaranteed to take x different tuples.

In opposition, with experience replay and a batch-size of 1, we are more likely to take the same tuple several times (less exploration): we don't extract as much from the trajectory. This effect is quite noticeable when we consider very few episodes that are very short. Indeed, the q-learning algorithm considers initial condition of zero and more elements stay to zero with small batch-size. Making plots would not have been too difficult but too time consuming (not asked anyway). So I'll only display the J's in figure 3.

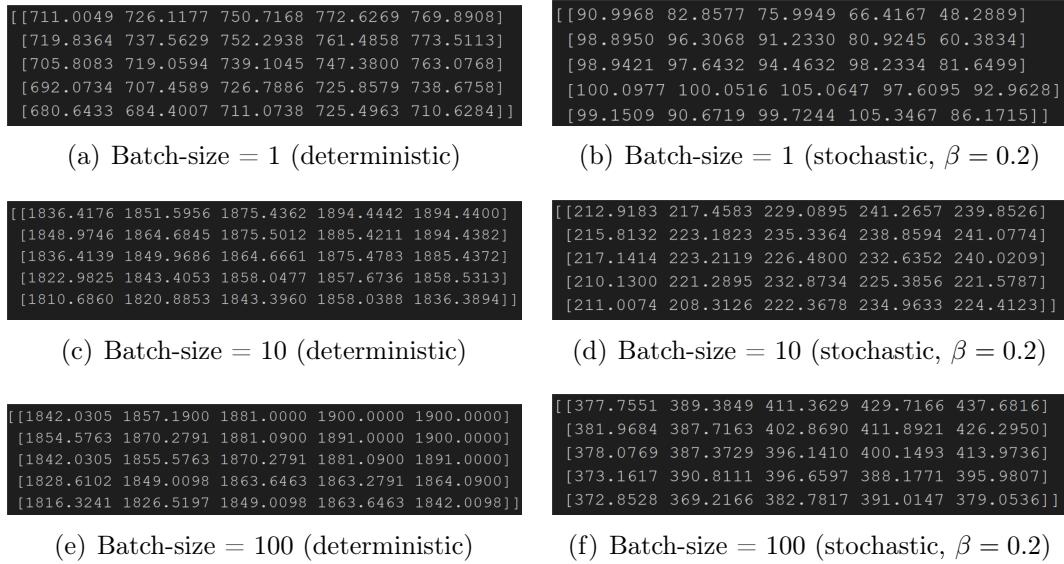


Figure 3 – Effect of batch-size on $J_{\hat{\mu}^*}^N$ when we consider experience replay. Number of samples = N is equal to 100000.

By playing around with more values than the ones exposed (1-10-100), it turns out that a small batch-size results in a very slow convergence. On the other hand, performance-wise, it is very fast to compute. Increasing a little bit the batch-size greatly benefit the speed of convergence and offer already good results with a batch-size of 10. However, this benefit is not much improved when considering significantly larger value of batch-sizes. It even turns out that large batch-sizes result in slow computations.

With regards to the distribution of the expected cumulative reward for each position in the grid, we can highlight the fact that, in the stochastic setting, the larger the batch-size, the larger is $J_{\hat{\mu}^*}^N(5, 5)$ compared to the other elements. For example, in figure 3(b), the expected cumulative result of $[5, 5]$ is the smallest value in the grid, while its immediate reward is 19=Br. By increasing the batch-size up to 100, we observe in figure 3(f) that it has reached the highest value in the grid, thanks to a more equitable sampling among the transitions present in the trajectory.

6.3 Analysis of effect of learning rate α

The protocol used is the one from the assignment with alpha as our parameter. The different \hat{Q} are displayed in figure 4 and figure 5 for deterministic and stochastic settings. In the same way, the expected cumulative rewards of the extracted optimal policy are displayed for each settings in figure 6.

[[705.506 723.913 716.937 717.396 706.976]	[[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[704.4005 724.4609 752.5300 777.9557 751.6551]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[724.4603 735.8510 754.6501 767.1332 779.5152]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[711.6651 714.4220 736.4500 747.4038 749.0614]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[699.6383 699.3624 726.2813 749.5126 723.5125]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[[703.2977 707.4529 722.9425 749.6182 775.7749]	[[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[724.4311 724.7863 737.1955 754.0674 767.6385]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[707.8929 712.2985 714.8924 737.3888 744.2071]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[699.5650 700.1274 695.0634 727.2821 747.9870]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[653.6226 657.3943 649.2536 701.1108 717.2407]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[[725.4078 744.6752 776.7869 797.6485 795.3585]	[[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[737.9435 754.6291 767.1325 779.4849 780.1579]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[714.7110 714.9214 746.9922 746.9280 750.6399]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[697.4494 724.7777 741.3806 721.9870 722.0488]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[666.1734 702.4634 718.7424 714.1864 714.5987]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[[726.0448 735.6479 755.6333 768.4777 783.4643]	[[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[712.9845 713.8638 738.4061 747.5442 751.0233]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[700.4232 698.7460 726.4665 745.9907 746.7792]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[660.3365 667.9105 704.7123 715.8988 713.1431]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]
[657.7813 666.8681 700.7463 716.3312 717.3344]	[1819.5722 1838.5004 1856.1497 1879.9529 1898.9399]

(a) $\hat{Q}, \alpha = 0.05$ (b) $\hat{Q}, \alpha = 0.2$ (c) $\hat{Q}, \alpha = 0.5$ (d) $\hat{Q}, \alpha = 0.5 \rightarrow 0$

Figure 4 – Effect of alpha on \hat{Q} in deterministic case. $[\uparrow, \leftarrow, \rightarrow, \downarrow]$ from top to bottom.

[[69.1673 57.3987 49.8591 39.4531 39.1876]	[[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[77.4777 67.3422 48.9966 49.8686 44.1073]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[97.8635 94.7682 84.1287 72.1494 50.6699]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[97.4845 83.2916 88.5174 81.8687 54.5117]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[100.6267 87.1780 93.9425 95.9294 75.8609]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[[71.6708 61.9536 53.3948 35.6882 31.7392]	[[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[90.7931 88.8416 75.6966 75.7397 53.8329]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[99.2100 94.3595 77.7733 81.5889 72.9984]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[99.2581 100.2392 85.0408 89.9002 92.7272]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[79.6386 78.9611 70.5664 83.7486 82.6785]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[[63.6587 47.4024 43.3808 44.9438 45.9775]	[[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[83.0371 82.9940 74.0557 54.3598 42.8111]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[82.7760 88.8699 82.6810 56.6173 49.1574]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[85.4564 96.6931 98.1020 78.4432 72.9675]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[71.1143 86.6881 86.3089 84.8809 79.5991]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[[87.8627 74.4050 67.0781 60.4386 35.8946]	[[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[97.8563 75.2286 83.1713 73.3465 44.3406]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[98.1342 85.2300 92.8337 89.9057 68.2565]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[79.2403 69.8706 85.3127 84.4421 76.1913]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]
[79.6202 71.2117 86.3248 85.2576 82.6363]	[1804.9922 1053.5458 1869.2455 1880.0407 1889.9563]

(a) $\hat{Q}, \alpha = 0.05$ (b) $\hat{Q}, \alpha = 0.2$ (c) $\hat{Q}, \alpha = 0.5$ (d) $\hat{Q}, \alpha = 0.5 \rightarrow 0$

Figure 5 – Effect of alpha on \hat{Q} in stochastic case ($\beta = 0.2$). $[\uparrow, \leftarrow, \rightarrow, \downarrow]$ from top to bottom.

[[259.3979 257.2639 259.5981 245.7624 245.1966]	[[1590.3972 574.9038 592.2166 569.1386 599.0444]
[260.5816 258.5361 261.1695 255.1471 255.9715]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[269.8725 273.7492 281.0055 281.8969 268.2495]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[266.3202 264.4442 271.9626 276.6034 267.1631]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[266.2638 258.9672 267.0557 294.1055 274.6513]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[[259.1871 251.7977 257.1870 257.8551 245.0861]	[[1590.3972 574.9038 592.2166 569.1386 599.0444]
[271.6167 266.5659 270.2284 280.6692 270.7364]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[269.2456 272.3880 267.3291 275.8907 281.1451]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[267.6694 269.7585 266.9684 263.6347 262.4477]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[244.5303 245.4050 255.0930 264.2565 265.0154]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[[258.1184 259.7608 251.1571 257.1973 249.8349]	[[1590.3972 574.9038 592.2166 569.1386 599.0444]
[267.8156 282.7065 284.7495 259.6695 263.5005]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[267.0637 277.2979 289.6634 262.0518 260.6238]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[261.1007 273.7098 282.6552 272.4389 266.3332]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[[256.8915 264.5185 271.1455 260.9546 255.9387]	[[1590.3972 574.9038 592.2166 569.1386 599.0444]
[272.9409 261.4561 281.9468 285.5983 252.4669]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[270.3283 259.1372 278.3962 294.3979 271.3593]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[245.3545 249.5503 265.2992 276.8923 279.1019]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[250.7252 267.7500 280.2915 279.1628 271.4027]	[1590.3972 574.9038 592.2166 569.1386 599.0444]

[[259.5913 261.8548 267.1340 257.5633 256.6252]	[[1590.3972 574.9038 592.2166 569.1386 599.0444]
[264.1260 273.1325 270.3218 263.2337 262.0100]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[268.5818 264.1619 267.3280 271.7492 266.4314]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[269.1594 272.0366 261.7882 272.4638 268.5287]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[269.4475 264.1411 271.8277 277.4485 265.2724]	[1590.3972 574.9038 592.2166 569.1386 599.0444]
[[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]	[[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]
[1853.2653 1868.3670 1879.7725 1889.6763 1898.6667]	[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]
[1864.7133 1854.2477 1868.9566 1879.7294 1889.6769]	[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]
[1827.2871 1867.3458 1852.4548 1869.7267 1881.0681]	[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]
[1815.0898 1863.7539 1847.6626 1862.2920 1894.0661]	[1590.6837 1855.8659 1879.6626 1898.6580 1898.6589]

(a) $J_{\mu^*}^N, \alpha = 0.05$, det. (b) $J_{\mu^*}^N, \alpha = 0.2$, det. (c) $J_{\mu^*}^N, \alpha = 0.5$, det. (d) $J_{\mu^*}^N$, det.

(e) $J_{\mu^*}^N, \alpha = 0.05$, sto. (f) $J_{\mu^*}^N, \alpha = 0.2$, sto. (g) $J_{\mu^*}^N, \alpha = 0.5$, sto. (h) $J_{\mu^*}^N$, sto.

Figure 6 – Effect of alpha on J . $\beta = 0.2$ for stochastic case.

The larger alpha, the larger the impact of the update in the q-learning. In a deterministic setting, increasing alpha helps the algorithm converge faster from what we can see in figure 4. Indeed an $\alpha = 0.05$ might be too conservative, thus resulting in a uselessly slow convergence.

In the stochastic setting, the value of alpha has a big impact on the final result. An alpha too high might lead to divergence. From the result obtained, a value of alpha close to 0.2 seems to offer acceptable results. It would have been perhaps interesting to plot the square error with $J_{\hat{\mu}^* N}$ by varying the value of alpha. Finally, there isn't much difference between a setting where we consider $\alpha_k = 0$, $\forall k$ and when we consider an alpha that decrease linearly towards 0.

6.4 Implement another exploration technique¹.

A widely used methods for balancing exploration/exploitation is a policy that selects at each time step a random action with a fixed probability, $0 \leq \epsilon \leq 1$, instead of selecting greedily one of the learned optimal actions with respect to the Q-function.

$$\pi(s) = \begin{cases} \text{random action from } U(s), & \text{if } \zeta < \epsilon \\ \operatorname{argmax}_{u \in U(s)} Q(s, u), & \text{Otherwise} \end{cases} \quad (1)$$

My code was no thought from the start to handle such kind of policies (the policy parameter for my simulation is a pointer to function). All the changes are not difficult but quite time consuming. For this reason, I only implemented it (`exploratoryPolicy()`) but didn't test it. Anyway, we can see that this policy offer a tradeoff between exploration and exploitation of (detected) high-reward paths. This tradeoff can be controlled via ϵ . In the previous stages, the trajectories were everytime generated by a random policy. Using (1) with $\epsilon < 1$ would encourage the agent to exploit more paths that are (currently) thought to give better rewards. This way, we might not discover some interesting cell with high rewards. It is more conservative in a sense.

Part II

$$\gamma = 0.4$$

With a smaller discount factor, the agent will favour closer but perhaps smaller rewards rather than long paths leading to a future excellent reward. In the same idea, the cumulative results will be smaller as new immediate rewards will rapidly vanish.

4 Expected return of a policy

One just need to change the variable GAMMA.

5 Optimal policy

Considering all the plots that have been made for each report, I'll just generate the more interesting results. If the full results were expected, one just needs to change `discntFact` in the main script to the new value and the remaining is automatic. The same applies for `plot.py`.

¹See *Adaptive ϵ -greedy Exploration in ReinforcementLearning Based on Value Differences*, Michel Tomic

```

Deterministic setting:
t=500
J:
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]

u:
[['down' 'down' 'down' 'right' 'up']
 ['left' 'right' 'right' 'right' 'up']
 ['up' 'right' 'up' 'down' 'up']
 ['left' 'right' 'right' 'left' 'left']
 ['up' 'right' 'up' 'up' 'right']]

Stochastic setting:
t=500
J:
[[1.3333 -0.0437 4.5302 16.6275 16.6275]
 [1.3333 4.5302 6.3235 9.4275 16.6275]
 [1.3333 0.7563 4.5302 12.3494 9.4275]
 [1.3333 3.2585 12.3494 3.2585 12.3494]
 [1.3333 -6.0506 3.2585 12.3494 4.8627]]]

u:
[['down' 'down' 'down' 'right' 'up']
 ['left' 'right' 'right' 'right' 'up']
 ['up' 'right' 'up' 'down' 'up']
 ['left' 'right' 'right' 'left' 'left']
 ['up' 'right' 'up' 'up' 'right']]

```

Figure 7 – μ^* and $J_{\mu^*}^N$ for each initial state x . $N=500$ in both deterministic and stochastic settings.

6 Q-learning in batch setting

Like expected, the convergence arise way sooner !

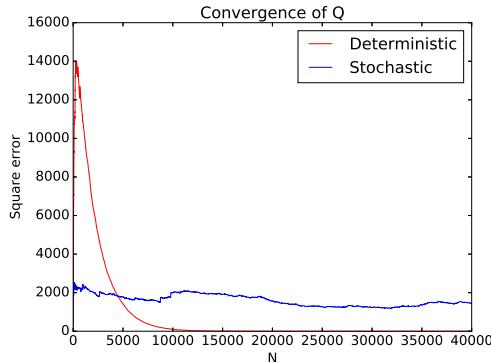


Figure 8 – Convergence of Q (equivalent of figure 1 with gamma = 0.4)

Question 6 with $\gamma = 0.4$

Q6: Q-learning in batch setting

1) Simulation of q-learning algorithm, 100 episodes of length 1000, starting in $[3, 3]$ by applying a random action at each step.

1.a) alpha = 0.05

1.a.i) Deterministic setting:

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'down']]
```

1.a.ii) Stochastic setting:

Expected cumulative reward of optimal policy:

```
[[2.6963 1.7177 2.3660 8.8219 12.7837]
 [1.3912 1.3954 6.0665 9.8261 13.8524]
 [4.2062 2.5416 5.3797 16.9032 10.0045]
 [3.6810 6.4285 16.5621 3.7236 11.9227]
 [-1.2166 -4.7157 5.2528 10.5159 6.9900]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'righ']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'down' 'down' 'up']
 ['up' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'down']]
```

1.b) alpha = 0.2

1.b.i) Deterministic setting:

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

1.b.ii) Stochastic setting:

Expected cumulative reward of optimal policy:

```
[[5.2661 3.0231 4.0239 13.6202 18.8060]
 [0.7489 3.1313 8.8356 10.7541 25.1118]
 [5.0085 4.8098 8.5829 20.7414 15.2833]
 [0.7045 7.0767 11.4094 4.0826 11.0682]
 [0.4795 -1.7601 7.0224 17.7679 3.9597]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'down' 'down']]
```

```
['left' 'righ' 'righ' 'left' 'up']
['up' 'left' 'up' 'down' 'up']
['up' 'righ' 'righ' 'left' 'left']
['up' 'righ' 'up' 'up' 'righ']]
```

1.a) alpha = 0.5

1.c.i) Deterministic setting:

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

1.c.ii) Stochastic setting:

Expected cumulative reward of optimal policy:

```
[[7.7330 6.4023 8.2341 30.8754 30.0656]
 [5.0079 9.2056 12.6448 14.7653 21.5975]
 [9.2723 6.5256 11.0820 22.8140 -1.4815]
 [1.1520 12.0694 22.8239 11.2308 19.4273]
 [3.4269 -1.5160 10.3606 22.7875 14.9768]]
```

Corresponding optimal policy:

```
[['down' 'down' 'righ' 'righ' 'righ']
 ['righ' 'righ' 'righ' 'righ' 'righ']
 ['down' 'righ' 'down' 'down' 'up']
 ['up' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

1.a) alpha decreasing linearly

1.c.i) Deterministic setting:

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
 [[1.0000 4.6363 1.0907 12.6667 31.6667]
 [1.0000 4.6363 1.0907 12.6667 31.6667]]
```

```

[10.0000 9.0907 15.2267 18.0667 22.6667]
[9.0000 -3.9637 10.0907 10.8095 1.0667]
[10.0000 -3.4762 13.8095 24.5238 4.8095]]

[[1.0000 1.0000 4.6363 1.0907 12.6667]
 [10.0000 10.0000 9.0907 15.2267 18.0667]
 [9.0000 9.0000 -3.9637 10.0907 10.8095]
 [10.0000 10.0000 -3.4762 13.8095 24.5238]
 [-16.0000 -16.0000 -16.3905 1.5238 6.8095]]

[[4.6363 1.0907 12.6667 31.6667 31.6667]
 [9.0907 15.2267 18.0667 22.6667 22.6667]
 [-3.9637 10.0907 10.8095 1.0667 1.0667]
 [-3.4762 13.8095 24.5238 4.8095 4.8095]
 [-16.3905 1.5238 6.8095 15.0000 15.0000]]

[[10.0000 9.0907 15.2267 18.0667 22.6667]
 [9.0000 -3.9637 10.0907 10.8095 1.0667]
 [10.0000 -3.4762 13.8095 24.5238 4.8095]
 [-16.0000 -16.3905 1.5238 6.8095 15.0000]
 [-16.0000 -16.3905 1.5238 6.8095 15.0000]]

```

1.c.ii) Stochastic setting:

Expected cumulative reward of optimal policy:

```

[[8.8947 2.7781 10.1388 25.2769 31.5251]
 [7.2838 9.4973 14.6396 22.0269 31.0267]
 [8.8095 7.3726 8.9653 14.5219 22.1129]
 [8.5267 4.2773 10.9394 9.6199 14.3541]
 [0.7482 -2.0071 0.5917 22.6793 14.4993]]

```

Corresponding optimal policy:

```

[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'down' 'down' 'up']
 ['up' 'righ' 'righ' 'left' 'down']
 ['up' 'righ' 'up' 'up' 'righ']]

```

Deterministic: Batch1

Expected cumulative reward of optimal policy:

```

[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]

```

Corresponding optimal policy:

```

[['down' 'down' 'down' 'righ' 'righ']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'down']]

```

Stochastic: Batch1

Expected cumulative reward of optimal policy:

```

[[3.7588 -1.6807 3.4400 15.7522 11.7797]
 [2.1960 3.4533 2.6733 12.9969 14.1840]
 [4.1975 0.5979 1.7023 12.3552 8.3686]
 [0.1077 3.6132 11.0404 1.3677 15.8848]
 [1.4705 -5.6387 2.4179 16.2660 7.0704]]

```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'righ']
 ['down' 'righ' 'righ' 'righ' 'up']
 ['down' 'left' 'down' 'down' 'up']
 ['up' 'righ' 'righ' 'up' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

Deterministic: Batch10

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

Stochastic: Batch10

Expected cumulative reward of optimal policy:

```
[[1.8172 -0.6248 2.2034 11.8197 17.0160]
 [2.1408 5.0508 2.1100 6.6810 9.6285]
 [2.7964 1.6744 5.1554 9.9312 6.2705]
 [4.4155 7.8763 16.2130 5.1266 11.1017]
 [6.2715 -5.2509 3.4604 9.7152 7.0726]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['righ' 'righ' 'righ' 'left' 'up']
 ['down' 'righ' 'up' 'down' 'up']
 ['up' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'down']]
```

Deterministic: Batch100

Expected cumulative reward of optimal policy:

```
[[10.0000 9.0907 15.2267 31.6667 31.6667]
 [10.0000 15.2267 18.0667 22.6667 31.6667]
 [10.0000 10.0907 15.2267 24.5238 22.6667]
 [10.0000 13.8095 24.5238 13.8095 24.5238]
 [10.0000 1.5238 13.8095 24.5238 15.0000]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'up']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['up' 'righ' 'up' 'down' 'up']
 ['left' 'righ' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'righ']]
```

Stochastic: Batch100

Expected cumulative reward of optimal policy:

```
[[2.4465 0.2912 4.7482 17.0682 21.8820]
 [3.4003 3.2955 5.4534 11.7766 22.9268]
 [4.0921 1.5500 4.5114 9.0486 11.3238]
 [1.7023 0.9362 14.0407 5.0345 12.1158]
 [0.8815 -6.0684 3.1835 16.3987 6.3874]]
```

Corresponding optimal policy:

```
[['down' 'down' 'down' 'righ' 'righ']
 ['left' 'righ' 'righ' 'righ' 'up']
 ['down' 'left' 'up' 'down' 'up']
 ['left' 'left' 'righ' 'left' 'left']
 ['up' 'righ' 'up' 'up' 'down']]
```
