

1.

(a) It produced negative reward at first and zero for the rest.

(b) It produced small sized negative rewards for the first few times and produced positive rewards for the rest in average about 100~150.

ex) first ten and last ten

-19.5	208.5
-13.0	93.5
-15.5	180.0
-12.5	72.0
-11.5	104.0
-6.5	173.5
-13.0	164.5
-4.5	70.5
-4.0	154.5
-9.0	123.5
...->	

(c) It produced negative rewards for most of the time.

ex) first ten and last ten

-36.5	-49.5
-36.5	-58.0
-30.0	-13.0
-36.5	-58.0
-22.0	-23.5
-25.0	-39.5
-26.5	-42.5
-16.0	-25.5
-31.0	-36.0
-22.5	-28.5
...->	

2

(a) It produced negative rewards for all the time however for average about -10

ex) first ten and last ten

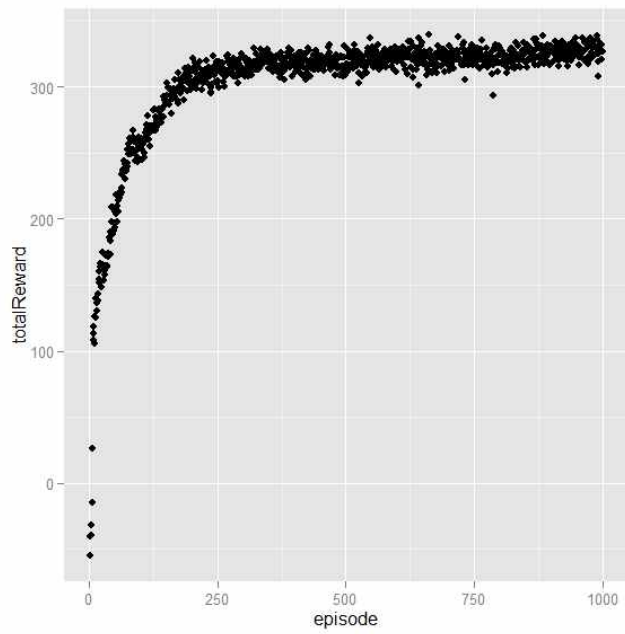
-40.5	-16.0
-24.5	-9.0
-17.0	-16.0
-12.0	-10.5
-19.0	-12.5
-11.5	-10.0
-9.5	-13.0
-14.5	-10.0
-12.5	-9.0
-13.0	-13.0
...->	

(b) It produced negative rewards for the first few and positive for the rest of simulation. The rewards increased in the beginning of the simulation until later 200's. Average about 200.

-71.0	275.5
-52.0	268.5
-50.0	278.5
-35.5	283.5
-23.5	275.5
11.0	279.5
61.0	279.5
83.0	297.5
82.0	295.0
102.5	299.0

(c) I have found best learning rate to be 0.1 and epsilon to be 0.01. This produces rewards about average 300.

3. The total reward increases each episode. Data is yet scattered a little.



For the averaged rewards, the graph looks similar to the first one, however, it looks neat since data are less scattered due to averaged process.

