# Assignment 1 - Defining  Solving RL Environments

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

This report aims to define and analyse a grid environment which follows OpenAI Gym standards. The environment can be configured to be deterministic or a stochastic one.

## 1 Checkpoint 1: Defining RL environments

### 1.1 Describe the deterministic and stochastic environments, which were defined (set of actions/states/rewards, main objective, etc).

The deterministic environment is a 4x4 grid(16 states) with (0,0) as the start point and goal at (3,3). The action space is of size 4, i.e., down, up, left, right. The max allowed time-steps is 10.

Our agent is a farmer with a donkey, the goal of the farmer is to reach the shed while loading all the bales of hay along the way on the donkey. The farmer wants to reach the shed quickly so he wants to avoid any patches of grass the donkey might start grazing.

There are 3 bales of hay/rewards in the grid(worth +1 each) and the 4th reward is the shed/goal itself(worth +3).There are 2 grass patches/obstacles in the grid(worth -1 reward each), further, after max time-steps have elapsed the environment will generate a reward of -3 for every subsequent action.

In the stochastic env, the only difference is that sometimes the donkey doesn't listen to the farmer and the requested action/step doesn't execute.
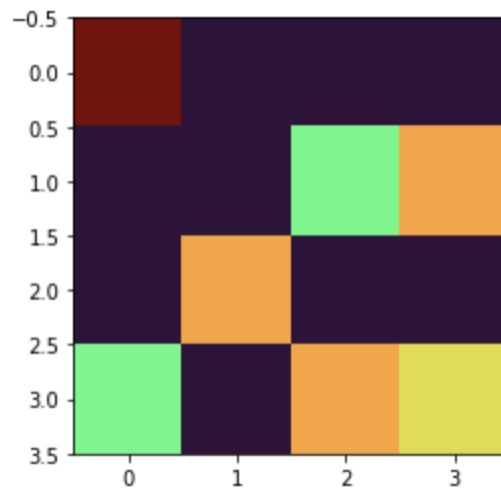
### 1.2 Provide visualizations of your environments.

Red - Agent
Green - Obstacle
Orange - Reward
Yellow - Goal

1

-0.5
0.0
0.5
1.0
1.5
2.0
2.5
3.0
3.5
     0     1     2     3

## 1.3    How did you define the stochastic environment?

The stochastic environment was defined by adding a 'stochastic' parameter to the environment. If this is set (self.stochastic = 1) there's 95% probability of the requested action being executed & a 5% probability that no action will be executed. The probabilities are based on what was recommended in the lectures.

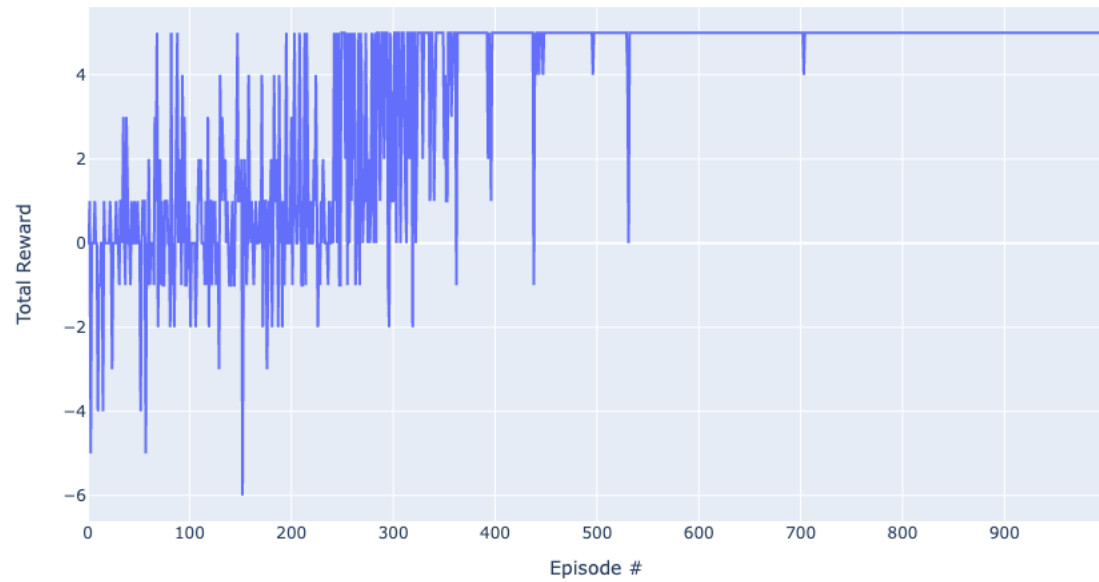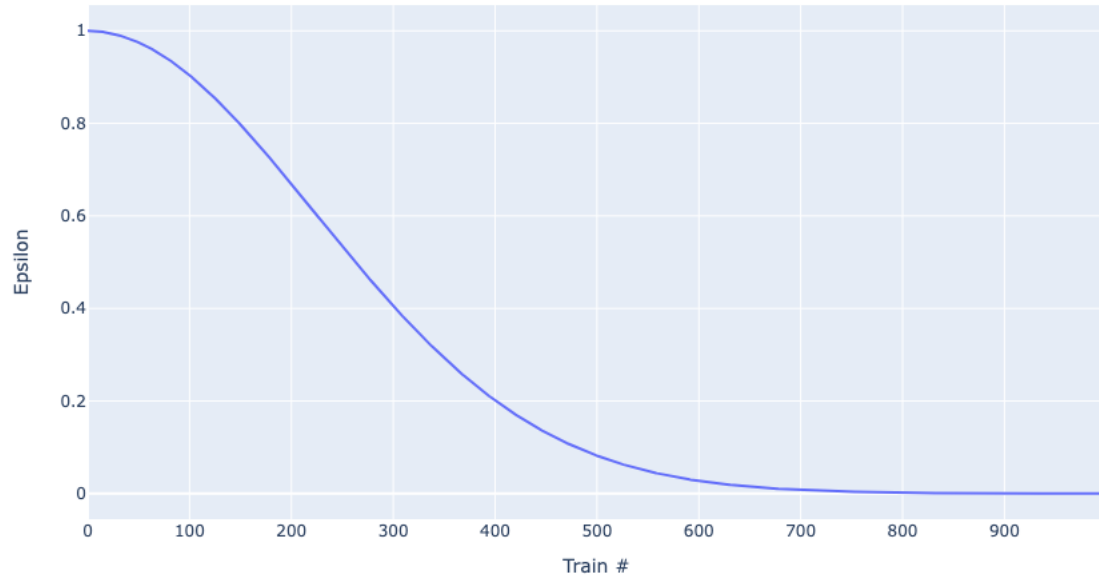## 1.4    What is the difference between the deterministic and stochastic environments?

In the stochastic env, sometimes the donkey doesn't listen to the farmer so a requested action doesn't execute and the farmer stays in the same state. If the farmer ends up remaining in the same state, he collects the reward of that state again.

## 1.5    Safety in AI: Write a brief review explaining how you ensure the safety of your environments.
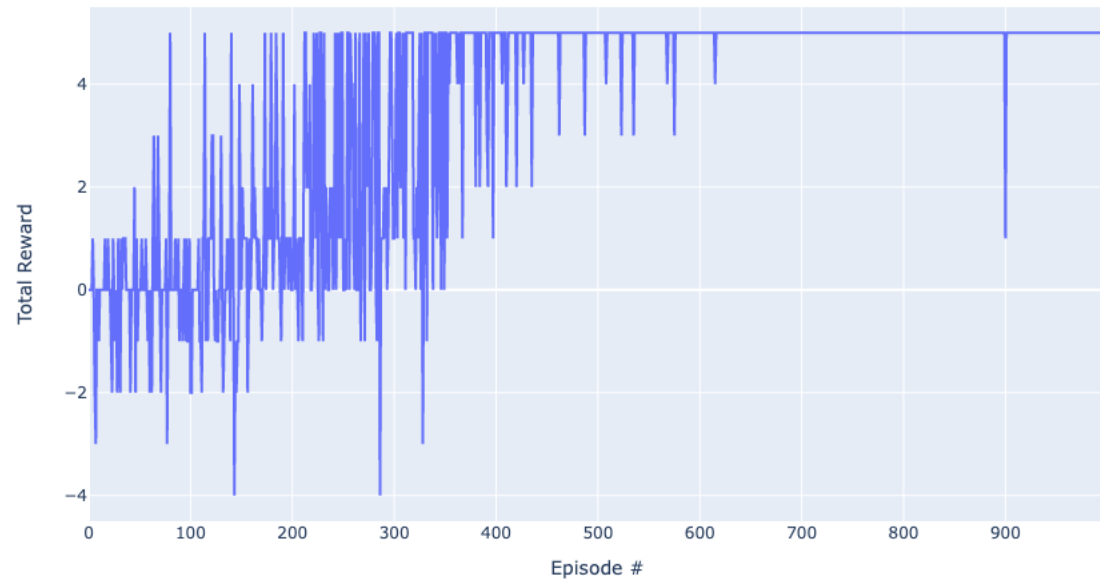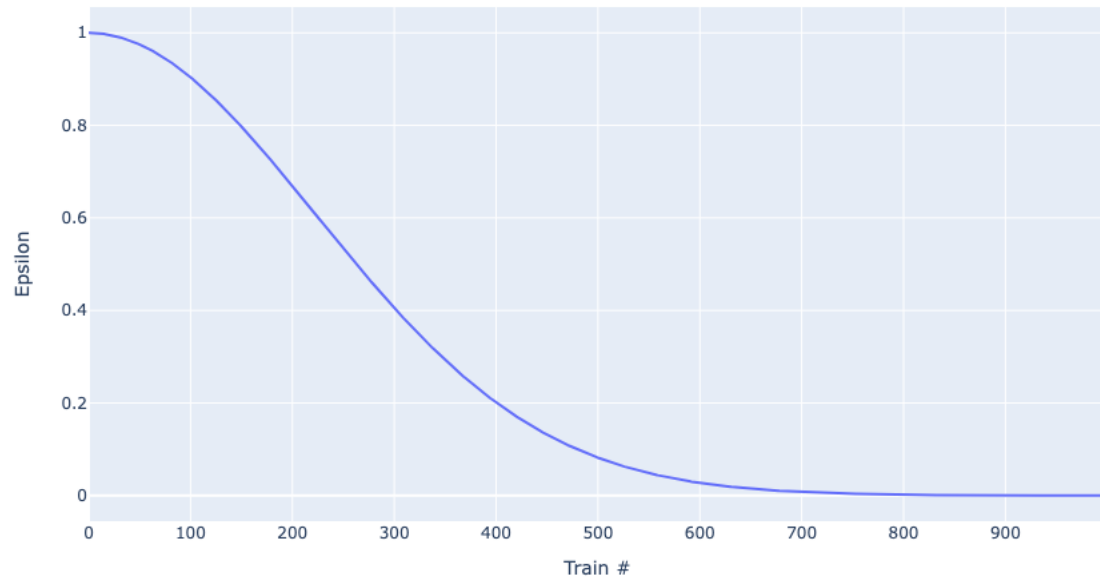
We ensure safety by using the np.clip() method. Using this method we clip the agent state list after every action such that the coordinates of the resultant position of the agent after taking an action always stays within the bounds of our environment.

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
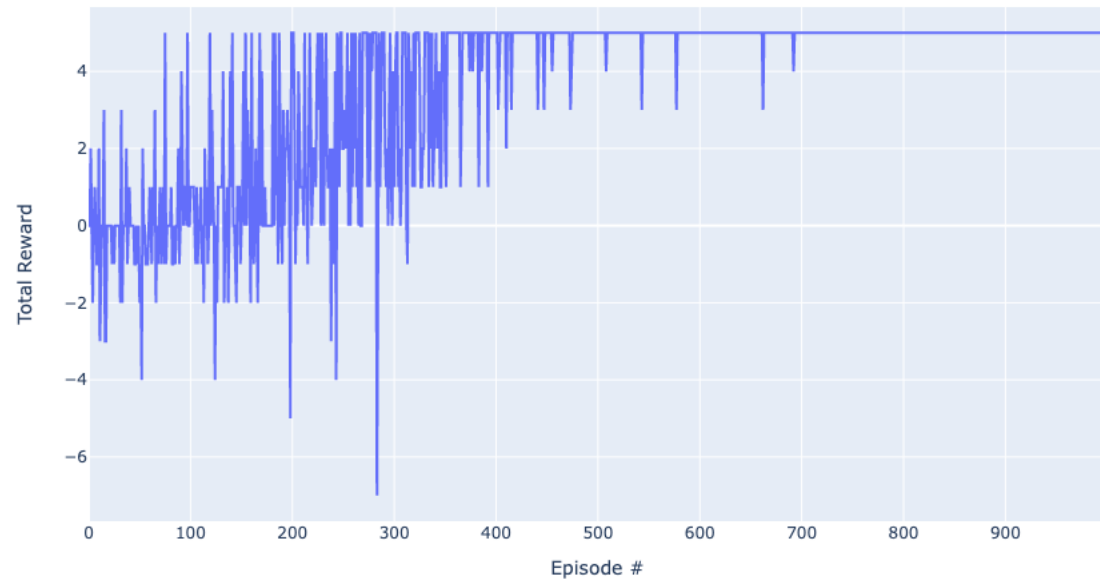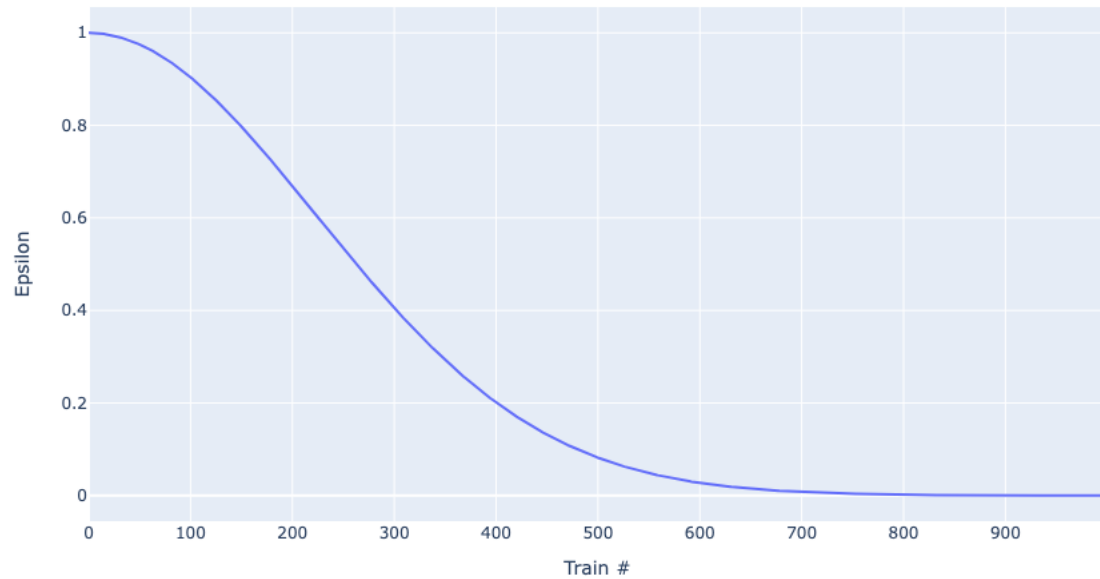160
161

# 2    Final submission: Tabular learning

## 2.1    Applying Q-learning to solve the deterministic environment defined in Part 1. Plots should include epsilon decay and total reward per episode.
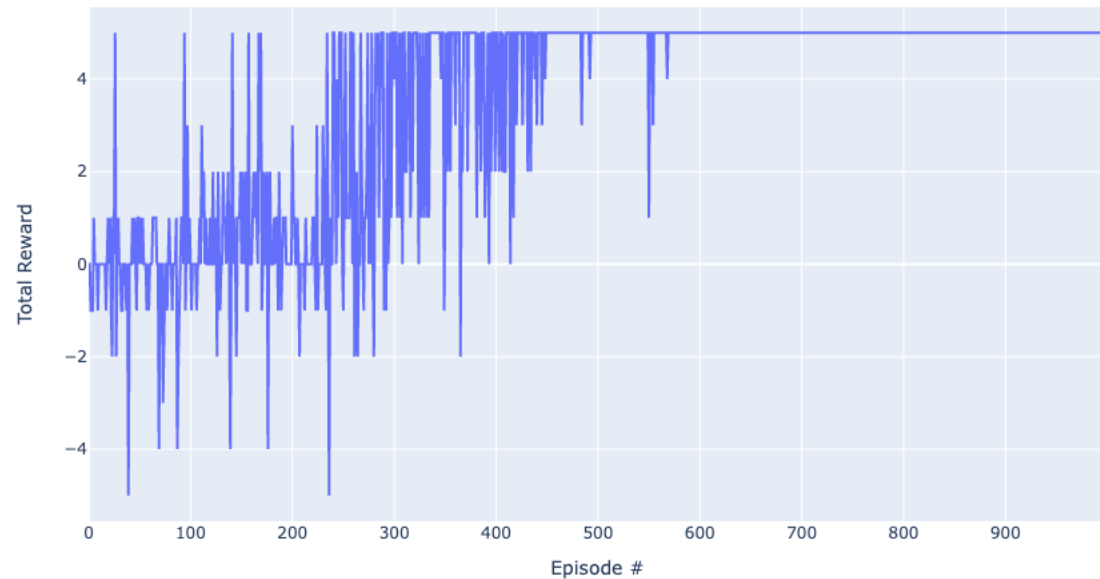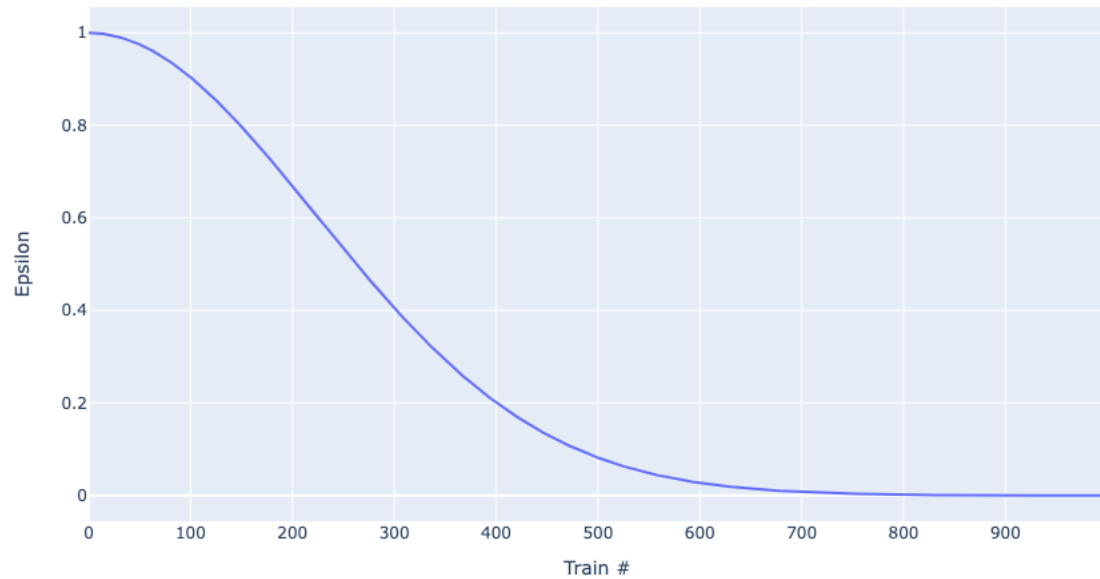
## 2.2 Applying Q-learning to solve the stochastic environment defined in Part 1. Plots should include epsilon decay and total reward per episode.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
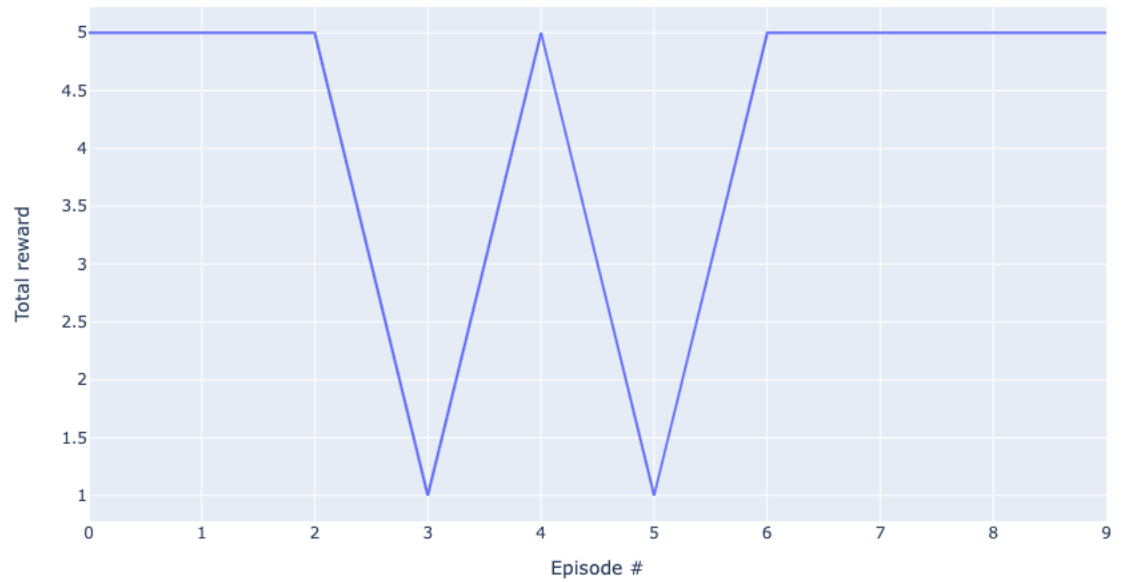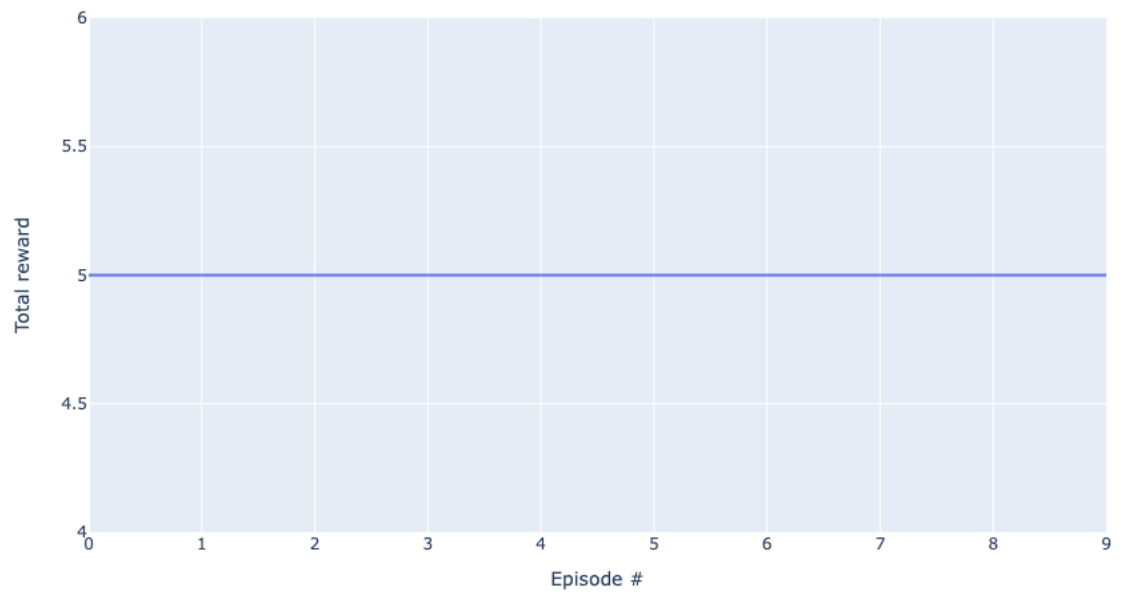259
260
261
262
263
264
265
266
267
268
269

## 2.3 Applying any other algorithm of your choice(SARSA) to solve the deterministic environment defined in Part 1. Plots should include total reward per episode.
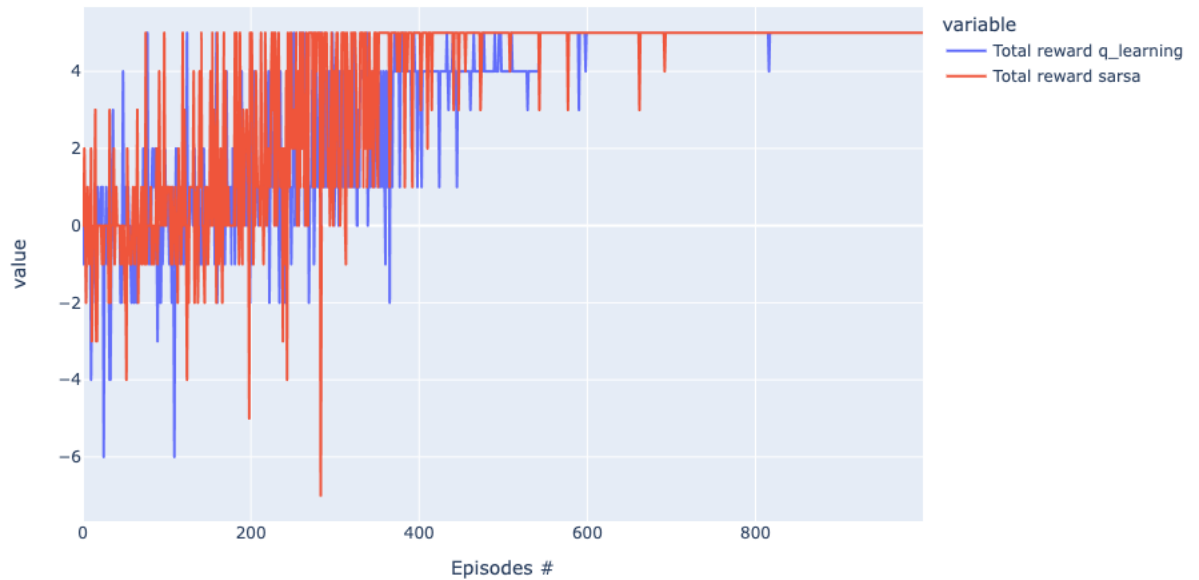
## 2.4 Applying any other algorithm of your choice(SARSA) to solve the stochastic environment defined in Part 1. Plots should include total reward per episode.
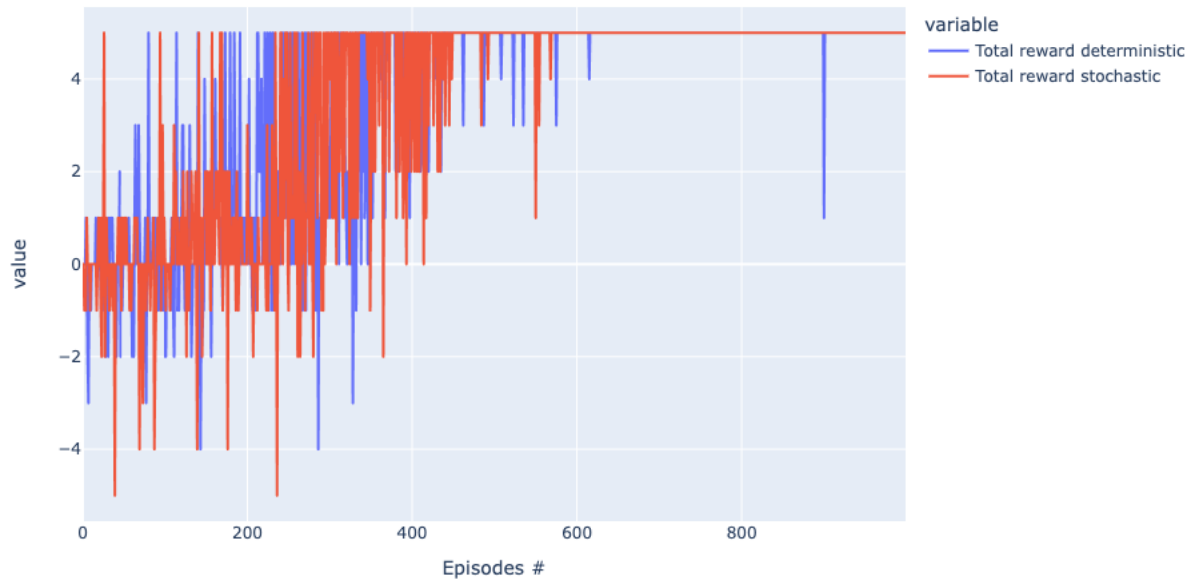
**2.5 Provide the evaluation results. Run your environment for at least 10 episodes, where the agent chooses only greedy actions from the learnt policy. Plot should include the total reward per episode.**

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

**2.6 Compare the performance of both algorithms on the same deterministic environment (e.g. show one graph with two reward dynamics) and give your interpretation of the results.**



**2.7 Compare how both algorithms perform in the same stochastic environment (e.g. show one graph with two reward dynamics) and give your interpretation of the results.**



8

## 2.8 Briefly explain the tabular methods, including Q-learning, that were used to solve the problems. Provide their update functions and key features.

### 2.8.1 Q-learning

In this off-policy tabular method, we use a greedy approach to converge towards the optimal policy.

Update func - qtable[state][action] = qtable[state][action] + alpha*(reward + discountfactor * maxq - qtable[state][action])

Here, maxq is the greedily chosen q value from the end state after we take an action.

### 2.8.2 SARSA

In this on-policy tabular method, we update the Q-table with policy determined decisions.

Update func - qtable[state][action] = qtable[state][action] + alpha*(reward + discountfactor * policyq - qtable[state][action])
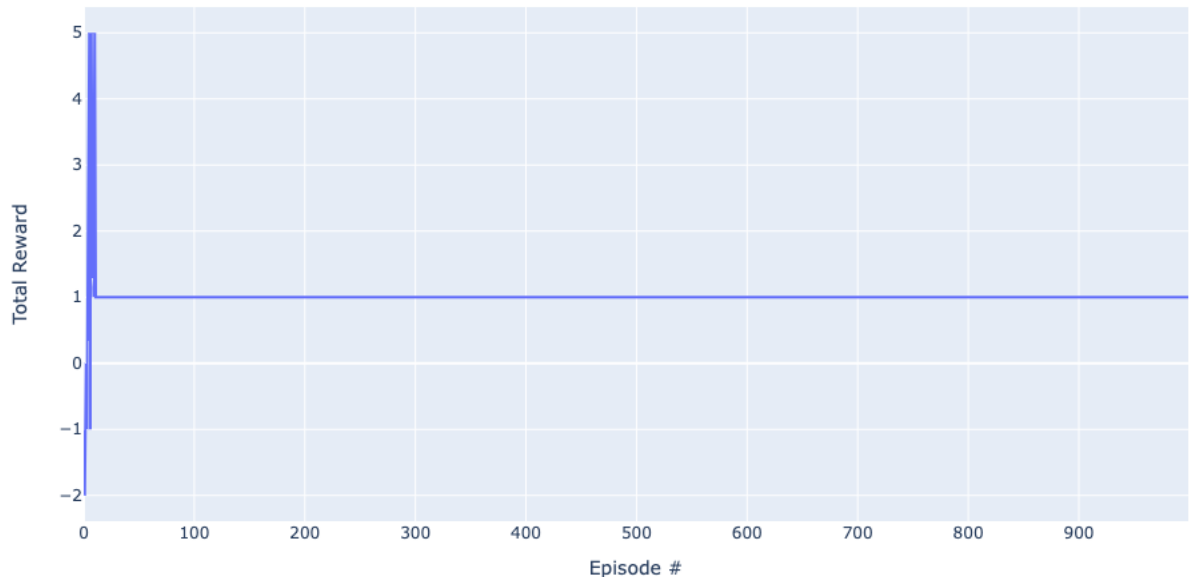
Here, the only difference from q-learning is that we choose 'policyq' based on the next action our policy(-greedy here) determines.

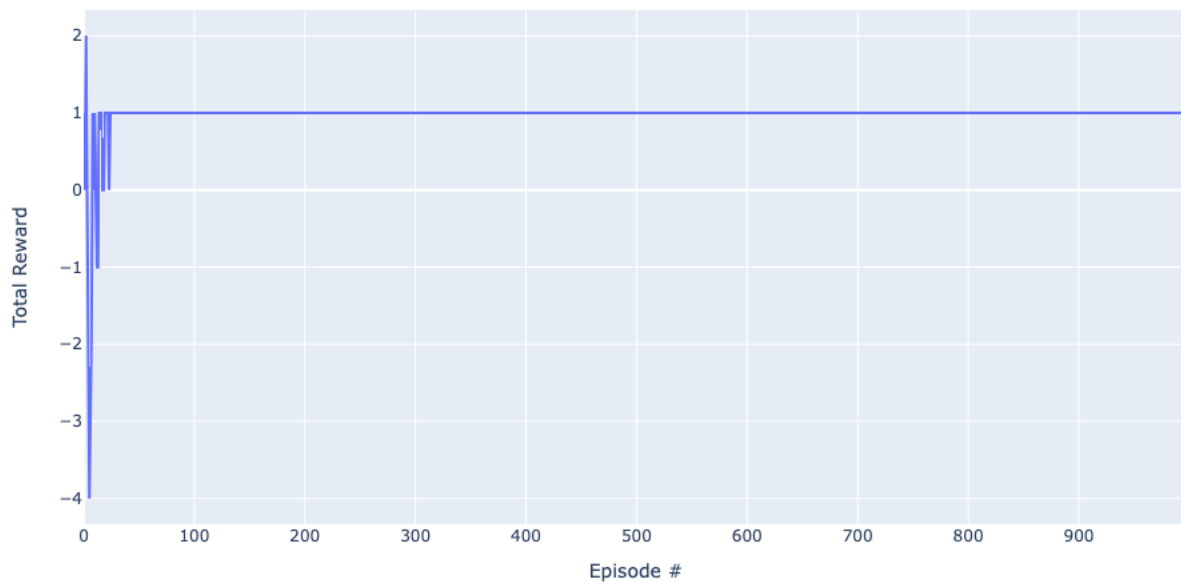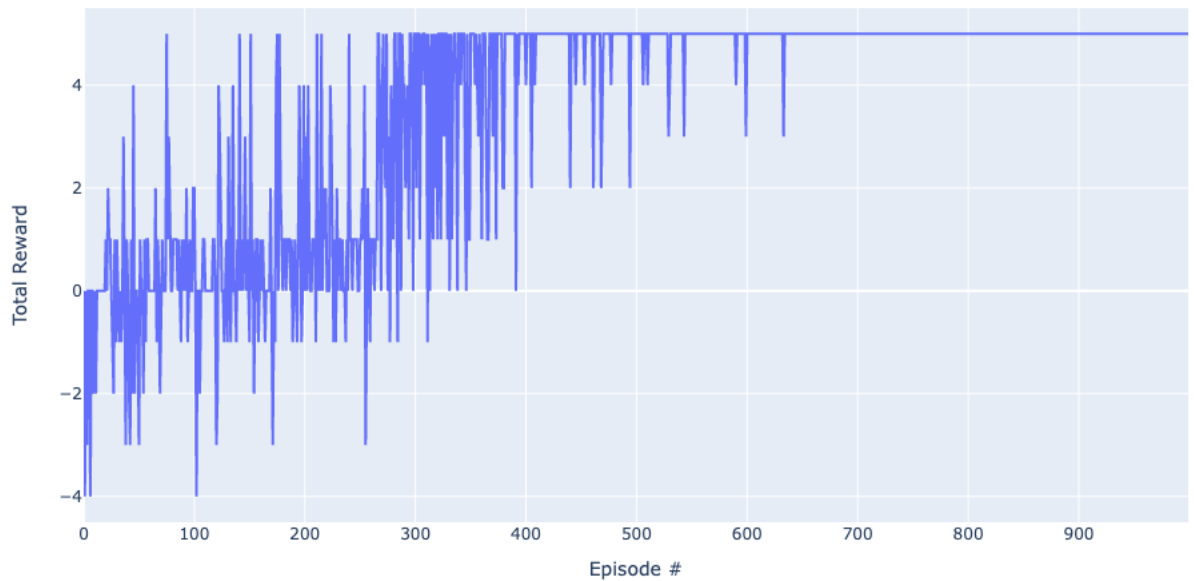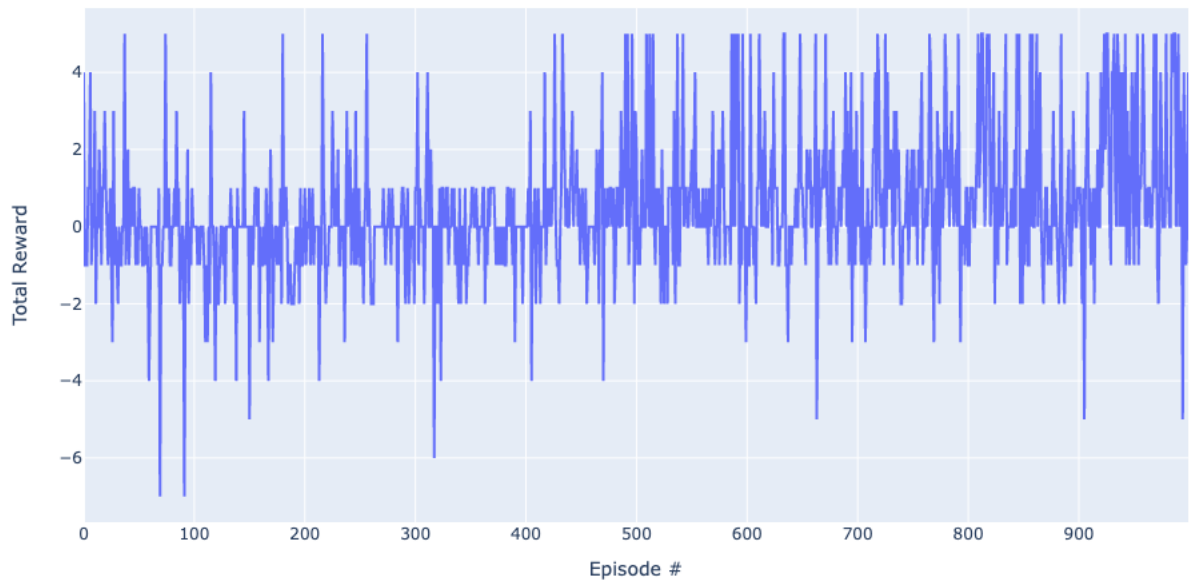# 3 Hyper parameter tuning

## 3.1 Epsilon decay rate

Please refer the jupyter notebook for more graphs, code etc.

When my epsilon was 0.03, my agent wasn't reaching the goal since the agent was exploring for long enough to find the best possible path as seen in the corresponding graph(steep slope). With the decay rate 0.000005, the slope was very gentle and my agent again wasn't reaching the goal as epsilon was still big enough towards the latter episodes to cause frequent spurts of random behaviour(exploration).

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
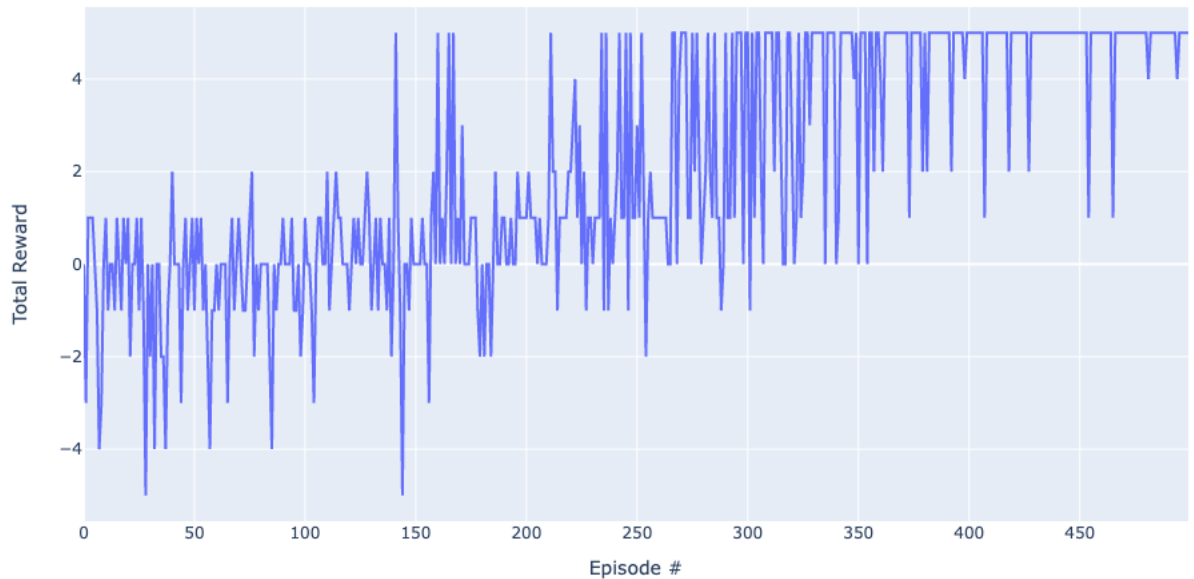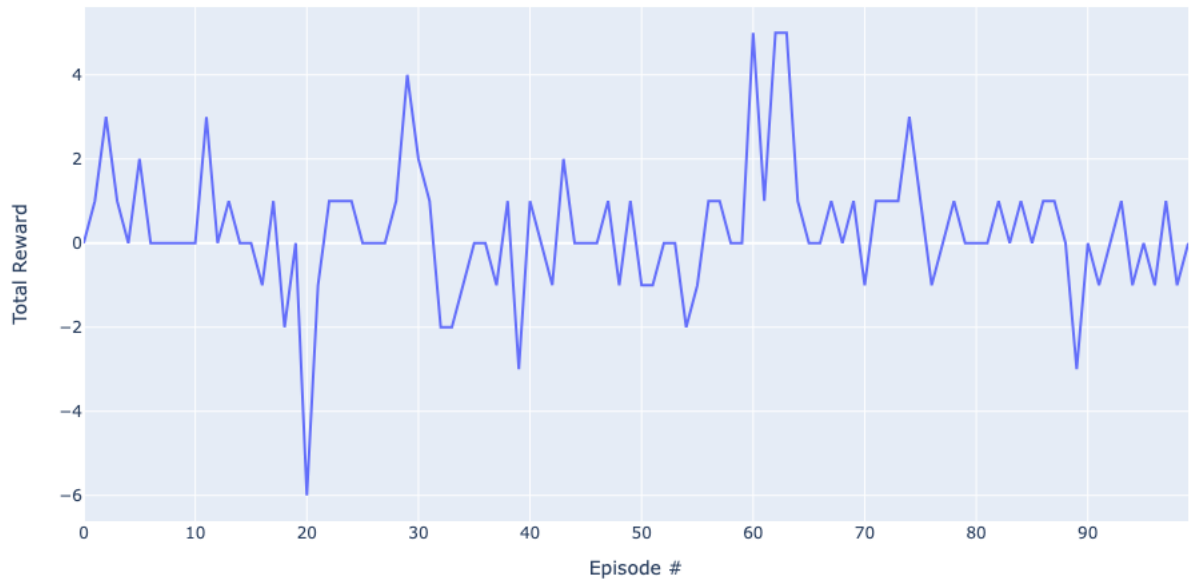579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
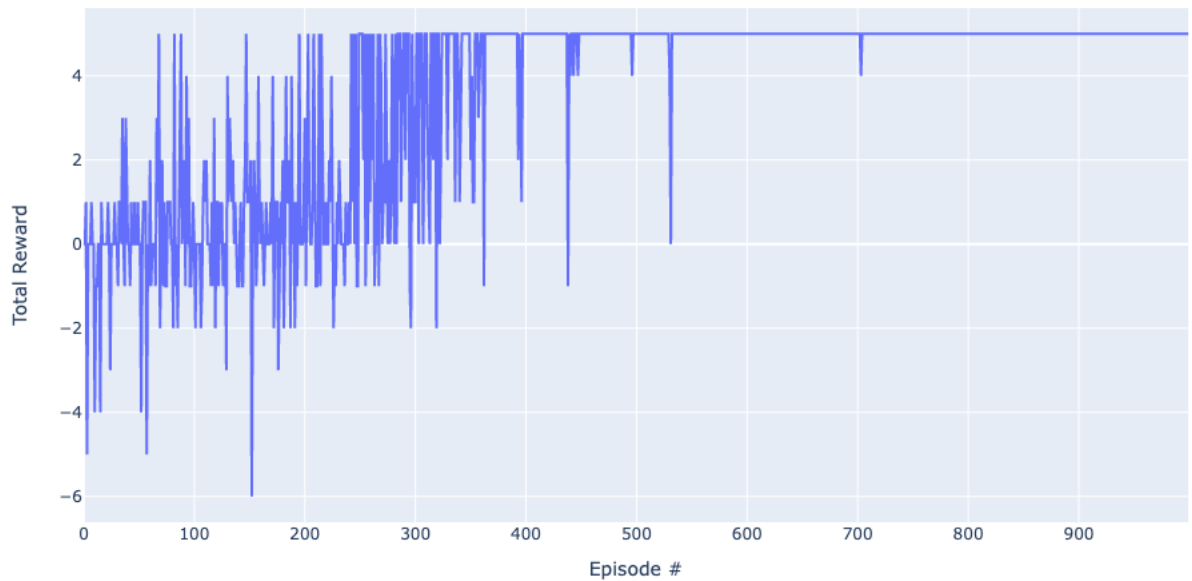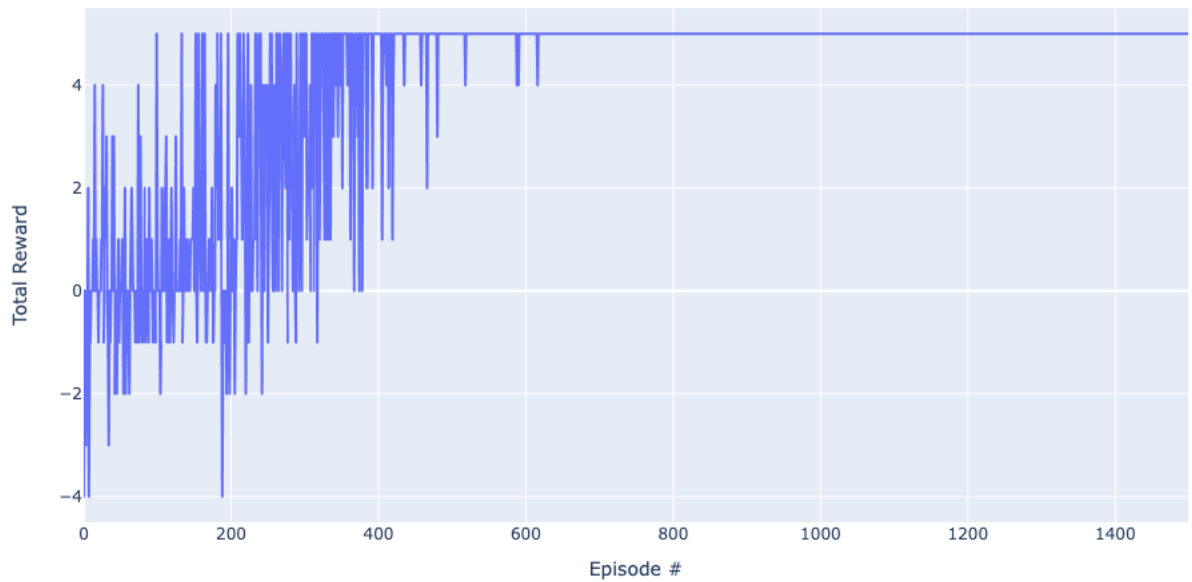
For epsilon=0.00002, the slope produced was much smoother causing a gradual transition to exploitation. Epsilon was converging to 0, roughly around 70% through the episodes( 700) which seemed like the right spot to only do exploitation. So this, was chosen as the decay rate for epsilon.

Note - I have used the exponential decay formula to calculate epsilon for each training episode.

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

## 3.2 Number of episodes

First we must note that number of episodes and epsilon decay are related in this case since we use exponential decay. For this experiment the epsilon decay rate was fixed to 0.00002 which was decided at optimum earlier.

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

When we train for 100 episodes we can see that rewards haven't convert much even in the latter iterations, for 500 episodes the rewards seem to have just converged towards the end episodes.

Third graph is 1500 iterations, where the rewards seem to converge around 600 mark so we just seem to be wasting iterations at this point. Finally for 1000, iteration we see that the rewards converge the same mark (a later fluctuation might be due to randomness in the epsilon greedy policy). Since that's also doesn't happen for this training after roughly 70

13