

# Machine Learning Nano Degree

## Project 4 – Train a SmartCab to drive

Daniel Fang 07/21/2016

### Basic Driving Agent

**QUESTION:** Observe what you see with the agent's behavior as it takes random actions. Does the *smartcab* eventually make it to the destination? Are there any other interesting observations to note?

When the agent only takes random actions, it seldom reached its final destination in the 100 trials. Besides, without taking rewards into consideration, it frequently breaks traffic rules. As shown in fig.1, the count of Succeed==True (which means the agent reached the destination in the trial) is about 10% of the total trials, and the median negative rewards is -11, which means the agent didn't obey the traffic rule at all.



Fig.1 Results of random actions

## Inform the Driving Agent

**QUESTION:** What states have you identified that are appropriate for modeling the *smartcab* and environment? Why do you believe each of these states to be appropriate for this problem?

The selected State model includes:

- 1) next\_waypoint : it allows the agent to keep awareness of its position relative to its final destination, which is the most important information to finish the task as soon as possible.
- 2) light : The lights tell us the state of the traffic light which is important that we add to state to ensure that our agent obeys the rule of the traffic light.
- 3) incoming & left & right : The agent needs to know the state of vehicles coming in on the left and on the right to help decide when to give way to traffic.

I didn't include deadline information for two reasons: 1) the time left to finish the task is not a very important feature compared to other features. The agent should not change its behaviors simply because there is a lot of time or limited time. 2) Including deadline information into state will significantly increase the possible number of states, which is too specific to the system and may not have good performance.

## Implement a Q-Learning Driving Agent

**QUESTION:** What changes do you notice in the agent's behavior when compared to the basic driving agent when random actions were always taken? Why is this behavior occurring?

Compared to random action agent, the Q-learning based agent is more intelligent in the way that its behavior became less erratic as the trials proceeded. Especially in the latter trials, the agent moved frequently towards its destination instead of moving the wrong directions. This is due to the Q-learning process, which 'teaches' the agent which action to take by selecting the highest rewards one.

## Improve the Q-Learning Driving Agent

**QUESTION:** Report the different values for the parameters tuned in your basic implementation of Q-Learning. For which set of parameters does the agent perform best? How well does the final driving agent perform?

The fine-tuning experiment is carried out on three variables: alpha, gamma. The results are shown in the following table. All experiment has an initial epsilon of 0.5, which decay asymptotically to 0 as trials proceed.

Experiment #	Alpha	Gamma	Success	Median_negative	Median_positive	Median_steps
1	0.9	0.9	96/100	0	22	13
2	0.9	0.4	97/100	0	22	13
3	0.9	0.2	100/100	0	22	13
4	0.4	0.9	48/100	-4	30	23
5	0.4	0.4	96/100	-0.5	24	13
6	0.4	0.2	98/100	0	22	13

From the above results, the optimal parameter set is the experiment #3.

**QUESTION:** Does your agent get close to finding an optimal policy, i.e. reach the destination in the minimum possible time, and not incur any penalties? How would you describe an optimal policy for this problem?

As result of this process we determined as the Optimal Policy the one obtained by using a Discount Factor (gamma) of 0.2 and a Learning Rate (alpha) of 0.9. With this policy, our Driving Agent successfully arrived to its final destination within the given timeframe in 100 out of 100 opportunities. The median number of actions is 13. Its median negative rewards is 0, but it still have some minor negative rewards. From the fig.2, we can see that the negative rewards has a decreasing tendency when trials proceeds, which has two reasons: first, the agents gained more experience with trials, and avoid negative rewards intentionally; Second, the exploration rate, epsilon, decays to almost zero in the latter part of the trials, which decrease the rate of agent making mistakes. However, there is still a chance that the agent may explore a little and that may incur penalties. And there are times that the agent did not follow the planer's suggested directions.

The optimal policy should be like this:

- The agent accumulates as much rewards as possible while still reaching its destination in time
- The agent should avoids negative rewards
- The agent should follows the planner's suggestion in most times.

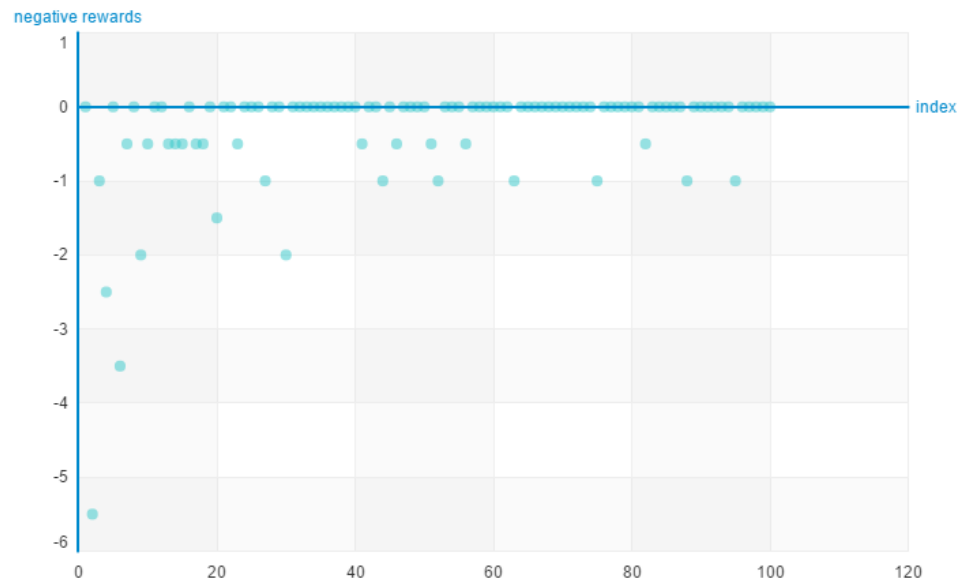


Fig.2 Negative rewards decreases as trials proceeds