

**Project 4**

# **Train a Smartcab to Drive**

Udacity Machine Learning Nanodegree Program

By Shinto Theruvil Manuel  
September 14, 2016

## Table of Contents

[Project Overview](#)

[Description](#)

[Software Requirements](#)

[Starting the Project](#)

[Running the Code](#)

[Definitions](#)

[Environment](#)

[Inputs and Outputs](#)

[Rewards and Goal](#)

[Performance Metrics](#)

[Task 1: Implement a Basic Driving Agent](#)

[Task 2: Inform the Driving Agent](#)

[Task 3: Implement a Q-Learning Driving Agent](#)

[Task 4: Improve the Q-Learning Driving Agent](#)

## Project Overview

In this project you will apply reinforcement learning techniques for a self-driving agent in a simplified world to aid it in effectively reaching its destinations in the allotted time. You will first investigate the environment the agent operates in by constructing a very basic driving implementation. Once your agent is successful at operating within the environment, you will then identify each possible state the agent can be in when considering such things as traffic lights and oncoming traffic at each intersection. With states identified, you will then implement a Q-Learning algorithm for the self-driving agent to guide the agent towards its destination within the allotted time. Finally, you will improve upon the Q-Learning algorithm to find the best configuration of learning and exploration factors to ensure the self-driving agent is reaching its destinations with consistently positive results.

## Description

In the not-so-distant future, taxicab companies across the United States no longer employ human drivers to operate their fleet of vehicles. Instead, the taxicabs are operated by self-driving agents — known as smartcabs — to transport people from one location to another within the cities those companies operate. In major metropolitan areas, such as Chicago, New York City, and San Francisco, an increasing number of people have come to rely on smartcabs to get to where they need to go as safely and efficiently as possible. Although smartcabs have become the transport of choice, concerns have arose that a self-driving agent might not be as safe or efficient as human drivers, particularly when considering city traffic lights and other vehicles. To alleviate these concerns, your task as an employee for a national taxicab company is to use reinforcement learning techniques to construct a demonstration of a smartcab operating in real-time to prove that both safety and efficiency can be achieved.

## Software Requirements

This project uses the following software and Python libraries:

- [Python 2.7](#)
- [NumPy](#)
- [PyGame](#)
  - Helpful links for installing PyGame:
  - [Getting Started](#)
  - [PyGame Information](#)
  - [Google Group](#)
  - [PyGame subreddit](#)

If you do not have Python installed yet, it is highly recommended that you install the [Anaconda](#) distribution of Python, which already has the above packages and more included. Make sure that you select the Python 2.7 installer and not the Python 3.x installer. `pygame` can then be installed using one of the following commands:

Mac:

```
conda install -c https://conda.anaconda.org/quasiben pygame
```

Windows & Linux:

```
conda install -c https://conda.anaconda.org/tlatorre pygame
```

## Starting the Project

For this assignment, you can find the `smartcab.zip` archive containing the necessary project files as a downloadable in the Resources section. *You may also visit our [Machine Learning projects GitHub](#) to have access to all of the projects available for this Nanodegree.*

This project contains two directories:

- `/images/`: This folder contains various images of cars to be used in the graphical user interface. You will not need to modify or create any files in this directory.
- `/smartcab/`: This folder contains the Python scripts that create the environment, graphical user interface, the simulation, and the agents. You will not need to modify or create any files in this directory except for `agent.py`.

In `/smartcab/` are the following four files:

- Modify:
  - `agent.py`: This is the main Python file where you will be performing your work on the project.
- Do not modify:
  - `environment.py`: This Python file will create the smartcab environment.
  - `planner.py`: This Python file creates a high-level planner for the agent to follow towards a set goal.
  - `simulation.py`: This Python file creates the simulation and graphical user interface.

## Running the Code

In a terminal or command window, navigate to the top-level project directory `smartcab/` (that contains the two project directories) and run one of the following commands:

```
python smartcab/agent.py or
```

```
python -m smartcab.agent
```

This will run the `agent.py` file and execute your implemented agent code into the environment. A README file has also been provided with the project files which may contain additional necessary information or instruction for the project. The following Definitions and Tasks slides will provide details for how the project will be completed.

## Definitions

### Environment

The smartcab operates in an ideal, grid-like city (similar to New York City), with roads going in the North-South and East-West directions. Other vehicles will certainly be present on the road, but there will be no pedestrians to be concerned with. At each intersection there is a traffic light that either allows traffic in the North-South direction or the East-West direction. U.S.

Right-of-Way rules apply:

- On a green light, a left turn is permitted if there is no oncoming traffic making a right turn or coming straight through the intersection.
- On a red light, a right turn is permitted if no oncoming traffic is approaching from your left through the intersection. To understand how to correctly yield to oncoming traffic when turning left, you may refer to [this official drivers' education video](#), or [this passionate exposition](#).

### Inputs and Outputs

Assume that the smartcab is assigned a route plan based on the passengers' starting location and destination. The route is split at each intersection into waypoints, and you may assume that the smartcab, at any instant, is at some intersection in the world. Therefore, the next waypoint to the destination, assuming the destination has not already been reached, is one intersection away in one direction (North, South, East, or West). The smartcab has only an egocentric view of the intersection it is at: It can determine the state of the traffic light for its direction of movement, and whether there is a vehicle at the intersection for each of the oncoming directions. For each action, the smartcab may either idle at the intersection, or drive to the next intersection to the left, right, or ahead of it. Finally, each trip has a time to reach the destination which decreases for each action taken (the passengers want to get there quickly). If the allotted time becomes zero before reaching the destination, the trip has failed.

## Rewards and Goal

The smartcab receives a reward for each successfully completed trip, and also receives a smaller reward for each action it executes successfully that obeys traffic rules. The smartcab receives a small penalty for any incorrect action, and a larger penalty for any action that violates traffic rules or causes an accident with another vehicle. Based on the rewards and penalties the smartcab receives, the self-driving agent implementation should learn an optimal policy for driving on the city roads while obeying traffic rules, avoiding accidents, and reaching passengers' destinations in the allotted time.

## Performance Metrics

Following table shows the three metrics identified for assessing the performance of our driving agent:

1	<b>Success Rate</b>	The primary objective of the driving agent is to reach the destination. This metric is the percentage of times the driving agent reaches the destination.
2	<b>Red light violations</b>	Plot of red light violations. Learning the rules of the road is an important part of learning to drive. This metric counts the serious mistake, red light violations and accidents that result in penalty of -1.0.
3	<b>Planner noncompliance</b>	Plot of planner noncompliance. Choosing an optimal route is what makes smartcab really smart. The planner will provide direction and then we will see whether smartcab will learn to follow the planner inputs. This metric counts the number times the driving agent not following the planner directions, that result in a penalty of -0.5.



## Task 1: Implement a Basic Driving Agent

To begin, your only task is to get the smartcab to move around in the environment. At this point, you will not be concerned with any sort of optimal driving policy. Note that the driving agent is given the following information at each intersection:

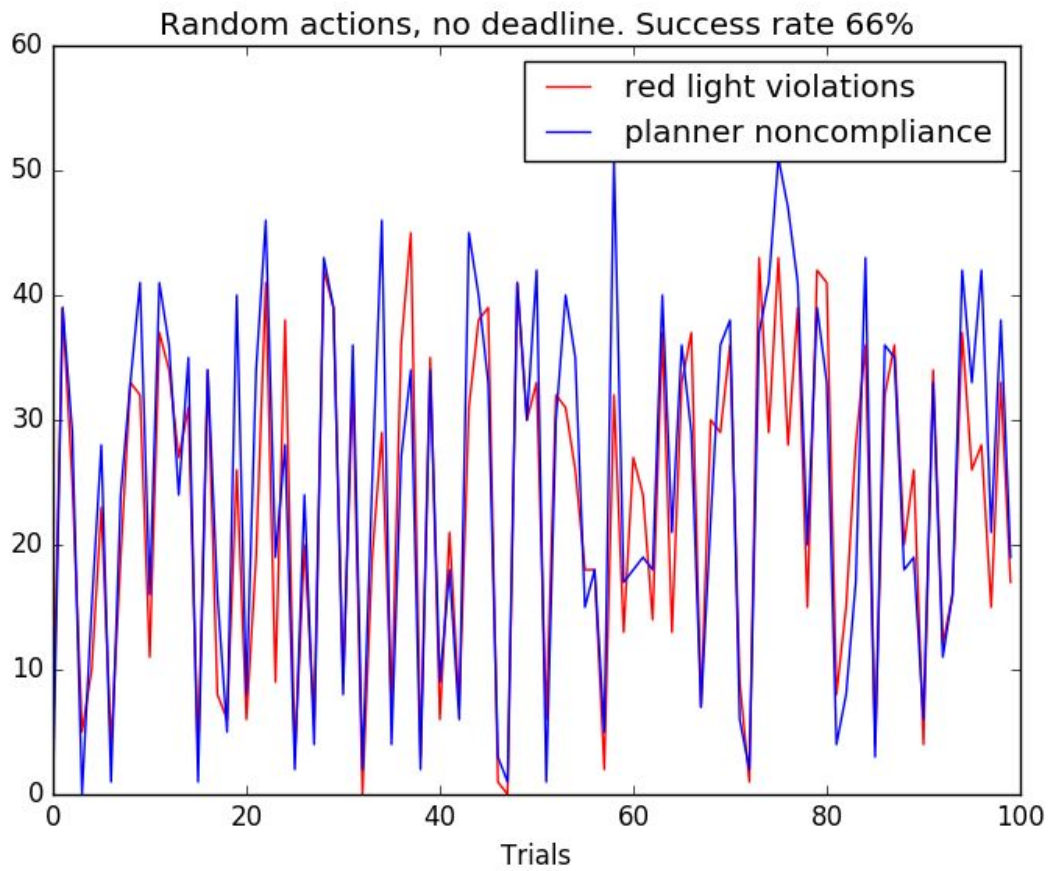
- The next waypoint location relative to its current location and heading.
- The state of the traffic light at the intersection and the presence of oncoming vehicles from other directions.
- The current time left from the allotted deadline.

To complete this task, simply have your driving agent choose a random action from the set of possible actions (**None**, **'forward'**, **'left'**, **'right'**) at each intersection, disregarding the input information above. Set the simulation deadline enforcement, **enforce\_deadline** to **False** and observe how it performs.

**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?*

The smartcab started moving. At each intersection the smartcab took one action randomly from four possible actions: None, forward, left, right. It was not considering environmental inputs. For example, many times it went through red light incurring penalties. On many trials the smartcab - reaches destination eventually.

The below plot of performance metrics shows that success rate as 66% without enforcing deadline. The plot has data on 100 trips the smartcab (or driving agent) has performed. There are two lines. The **red line** shows the **'red light violations'** in each trip. We can see that driving agent continue to make red light violations irrespective of the number trips it already performed. The driving agent is not learning from the experience. The **blue line** shows the **'planner noncompliance'**. The planner suggestions ( next waypoint) are like GPS guidance. Here also the driving agent is incurring penalties by not following the planner suggestions (next waypoint) throughout the 100 trips. The agent is not learning from the penalties or rewards.



If the deadline is enforced here, the success rate further reduces to an average of 16%.

## Task 2: Inform the Driving Agent

Now that your driving agent is capable of moving around in the environment, your next task is to identify a set of states that are appropriate for modeling the smartcab and environment. The main source of state variables are the current inputs at the intersection, but not all may require representation. You may choose to explicitly define states, or use some combination of inputs as an implicit state. At each time step, process the inputs and update the agent's current state using the `self.state` variable. Continue with the simulation deadline enforcement `enforce_deadline` being set to `False`, and observe how your driving agent now reports the change in state as the simulation progresses.

**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 agent (smartcab) can observe the following inputs from the environment::

- **Next waypoint:** This input is like GPS navigation. It helps the agent choose the next action towards the goal. If there were no constraints this would be the optimal action.
- **Light:** Traffic light. Possible values are Red, and Green
- **Oncoming:** Intention of traffic on the oncoming side of intersection.
- **Right:** Intention of traffic on the right side of intersection.
- **Left:** Intention of traffic on the left side of intersection
- **Deadline:** Number of available actions (steps) remaining.

### States identified as appropriate:

The following table shows the states identified as appropriate:

<b>Next way point</b>	Important in choosing the next action which leads towards the destination. Possible values are: None, “forward”, ‘right’, and ‘left’.
<b>Light</b>	Observing traffic light is important to avoid penalties.
<b>Oncoming</b>	Understanding the intention of oncoming traffic is important in taking a left turn on a green light, or taking right turn in a red light.
<b>Left</b>	Observing from left is important while taking a right turn on a red light.

Following two inputs are ignored at this stage:

**Right:** The intention of traffic on the right of intersection is not going to influence the smartcab decisions. So this input is seen redundant and so ignored.

**Deadline:** This input shows the concept of time. But not helpful in avoiding penalties or finding direction. This input is ignored as increases the number of states too much.

**OPTIONAL:** *How many states in total exist for the smartcab in this environment? Does this number seem reasonable given that the goal of Q-Learning is to learn and make informed decisions about each state? Why or why not?*

The states can be represented in four variables (next way point, light, oncoming and left) with the following possible values:

next way point  $\in$  { None, forward, left, right}  
 light  $\in$  { green, red}  
 oncoming  $\in$  { None, forward, left, right}  
 left  $\in$  { None, forward, left, right}

The size of the states become  $4 * 2 * 4 * 4 = 128$ . With the four actions for each state as:

The number of states, 128 is manageable. The Q-Learning algorithm requires only looking up and setting values in a matrix.

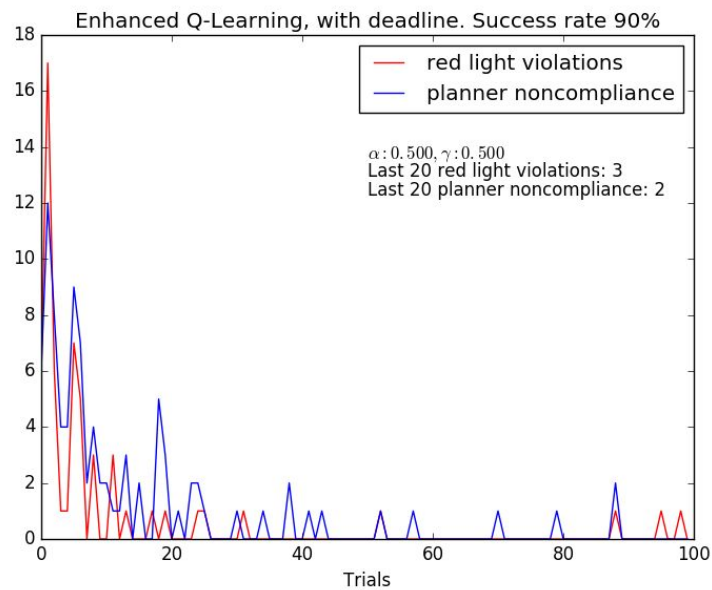
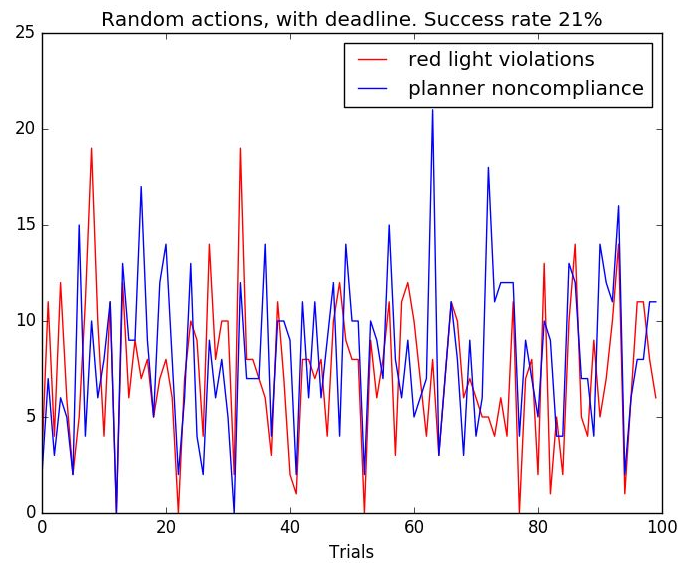
### Task 3: Implement a Q-Learning Driving Agent

With your driving agent being capable of interpreting the input information and having a mapping of environmental states, your next task is to implement the Q-Learning algorithm for your driving agent to choose the *best* action at each time step, based on the Q-values for the current state and action. Each action taken by the smartcab will produce a reward which depends on the state of the environment. The Q-Learning driving agent will need to consider these rewards when updating the Q-values. Once implemented, set the simulation deadline enforcement `enforce_deadline` to `True`. Run the simulation and observe how the smartcab moves about the environment in each trial.

The formulas for updating Q-values can be found in [this](#) video.

**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?*

From the two plots in next page, we can see the improvement in the 3 performance metrics. The success rate gone up close to 90 %. There is a clear reduction in red light violations. The agent is learning to take responsible actions if the light is red. The count of noncompliance to planner is also reduced as the number of trips progressed.



The above improvement is the result of Q-learning algorithm that guides the agent to select the best action.

Here is the steps of Q-Learning Algorithm:

1. Initialize the Q-table for the value,  $Q_j(s, a)$ , to 0s.
2. Choose the best action for the current state based on policy
3. Take action, observe the reward and new state

4. Update Q-value for the state using the observed reward and the maximum reward possible for the next state
5. Set the state to the state to new state
6. If not reached destination, repeat the steps

Q-Learning equation for updating the Q-value:

$$Q(s, a) = Q(s, a) + \alpha * [reward + \gamma * \max_a Q(s', a) - Q(s, a)]$$

**$\alpha$  (alpha)** - the learning rate, is a value between 0 and 1. Setting it to 0 means the Q-values are never updated, hence nothing is learned. Setting it to a high value such as 0.9 means that learning can happen fast but the algorithm might not converge. This parameter is set to 0.5

**$\gamma$  (gamma)** - the discount factor, is a value between 0 and 1. It control the value placed on future rewards. If the value is low, immediate rewards are prioritized , while higher values make it strive for long-term high reward. This parameter is set to 0.5

**$\epsilon$  (epsilon)** - the exploration factor is also a value between 0 and 1. Initially set to 0.9 and made to decay with a factor of 0.9 with each query in Q table.

With this basic Q-Learning implementation, 10 tests were done. Each test consists of 100 trials.

The performance metrics collected are:

- 1) **Success Rate** - Number of times the agent reached destination out of 100 trials
- 2) **Last 20 Red Light violations** - As the trials progress, the agent should learn to minimize penalties by avoiding red light violations. Here we count the red light violations in trials 80 to 100. This number should be close to zero for a perfectly learned agent.
- 3) **Last 20 Planner Noncompliances** - It shows towards end of trials whether the agent learned to follow follow directions from the planner. The agent should learn to minimize this value as the learning progresses. In a perfectly learned agent these value shold be close to zero

Here is the result of testing the baseline implementation of Q-Learning algorithm:

Test #	Success Rate	Red Light Violations in last 20 trials	Planner noncompliances in last 20 trials
1	90	4	67
2	98	3	11
3	88	3	97
4	92	1	7
5	95	3	9
6	93	0	9
7	83	2	114
8	69	7	48
9	90	3	2
10	90	1	91

#### Summary of Baseline :

	Success Rate	Last 20 Red Light Violations	Last 20 Planner noncompliances
<b>Mean</b>	88.8	2.7	45.5
<b>Minimum</b>	69	0	2
<b>Max</b>	98	7	114

#### Observations:

The agent has learned to follow the rules of the road. This is evident in the sharp reduction in the red light violations. The average red light violations recorded in the last 20 trials is just 2.7. On average 12.2 times the agent could not reach destination due to getting stuck in the local minima (agent not taking any action) or going in circles. The initial trails (left part of the plots) the agent incurs penalties, due to random exploration and later it learns to minimize it by selecting the best q values. We can also see that the blue line (planner non compliances) in the plot comes down as the trial progresses. This shows that the agent learns to use the 'next



waypoint' suggestions from planner.

## Task 4: Improve the Q-Learning Driving Agent

Your final task for this project is to enhance your driving agent so that, after sufficient training, the smartcab is able to reach the destination within the allotted time safely and efficiently.

Parameters in the Q-Learning algorithm, such as the learning rate (**alpha**), the discount factor (**gamma**) and the exploration rate (**epsilon**) all contribute to the driving agent's ability to learn the best action for each state. To improve on the success of your smartcab:

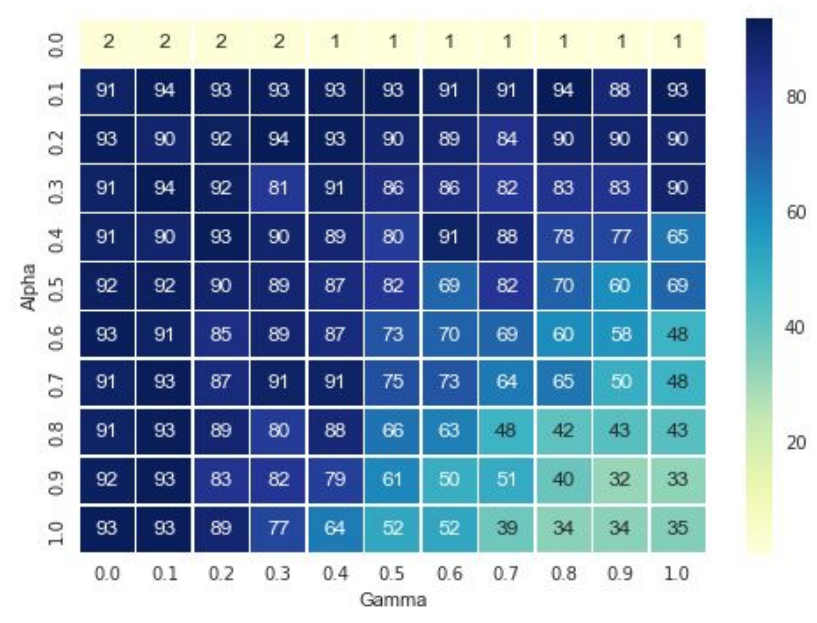
- Set the number of trials, **n\_trials**, in the simulation to 100.
- Run the simulation with the deadline enforcement **enforce\_deadline** set to **True** (you will need to reduce the update delay **update\_delay** and set the **display** to **False**).
- Observe the driving agent's learning and smartcab's success rate, particularly during the later trials.
- Adjust one or several of the above parameters and iterate this process.

This task is complete once you have arrived at what you determine is the best combination of parameters required for your driving agent to learn successfully.

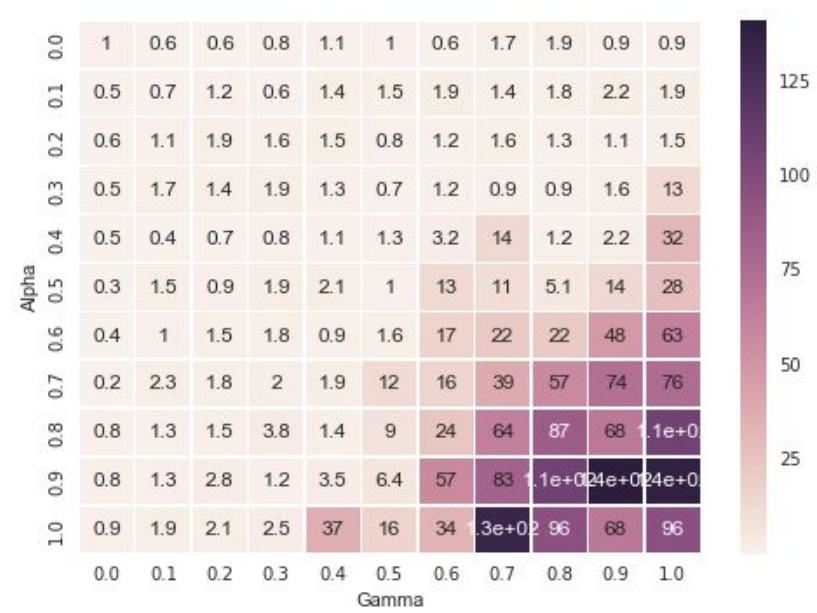
**QUESTION:** Report the different values for the parameters tuned in your basic implementation of Q-Learning. For which set of parameters does the agent perform best? How well does the final driving agent perform?

The Alpha and Gamma will take value between 0 and 1. To find the best value grid search is performed and heatmaps are prepared for better interpretation.

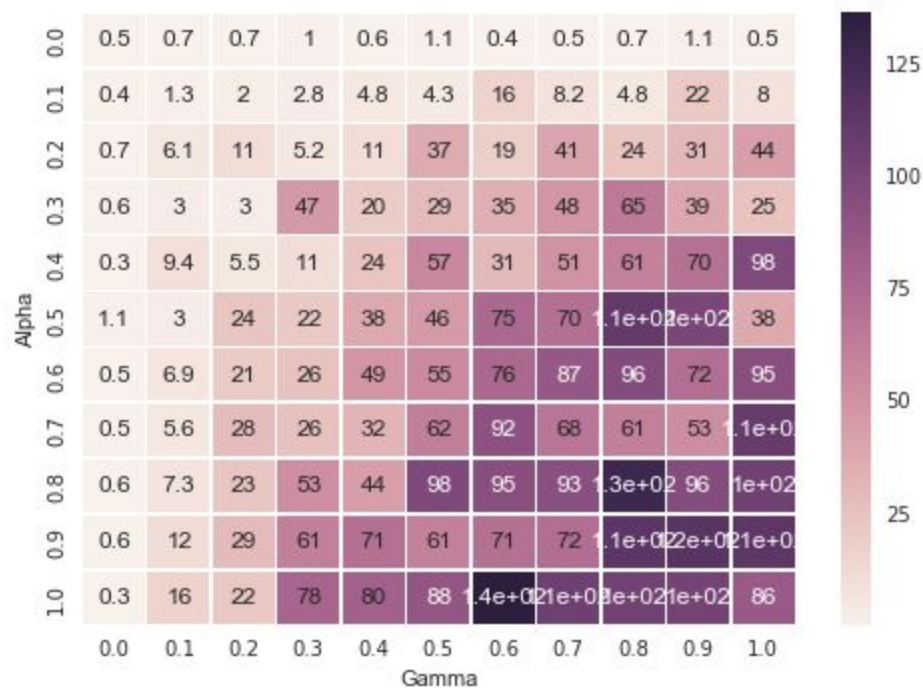
**Heatmap on Success Rate (Alpha & Gamma between 0 to 1)**



Heatmap on Red Light Violations (Alpha & Gamma between 0 to 1)

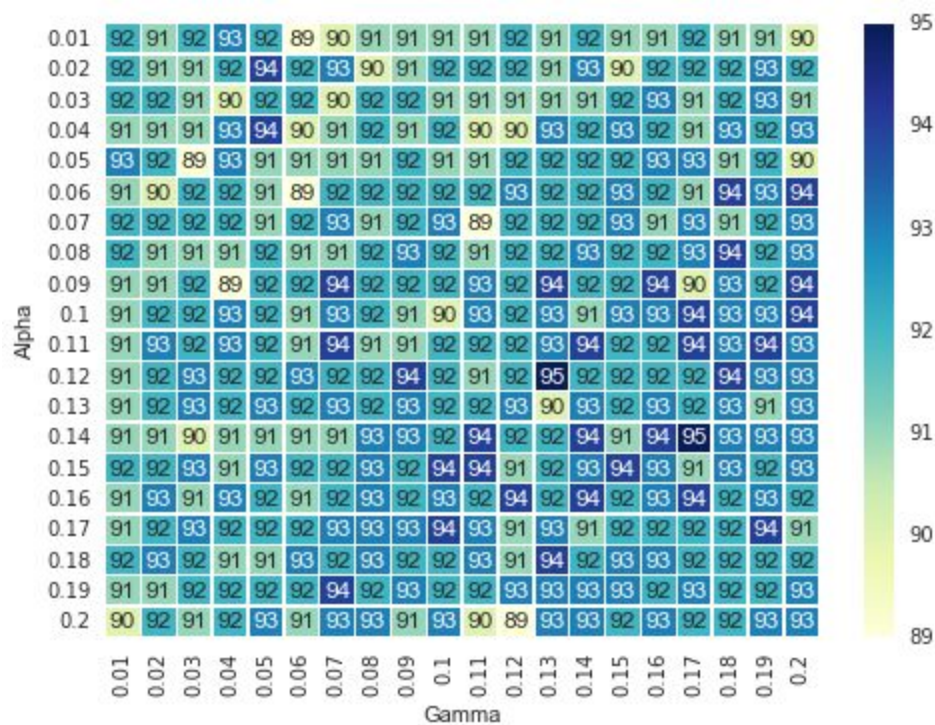


Heatmap of Planner Noncompliance (Alpha & Gamma between 0 to 1)

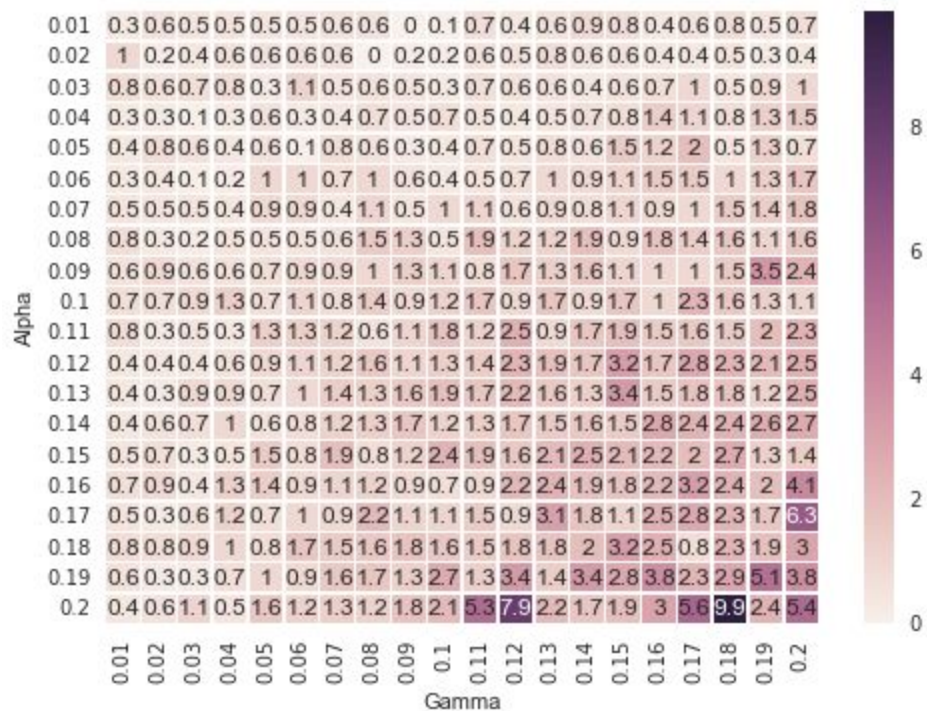


It looks like best results are towards the upper left corner, most likely a .value between 0.1 and 0.2. To narrow down further another grid search is done between 0.1 and 0.2 and heatmaps are generated.

#### Heatmap on Success rate (Alpha & Gamma between 0.01 to 0.2)



Heatmap on red light violations (Alpha & Gamma between 0.01 to 0.2)





Heatmap on Planner noncompliance (Alpha &amp; Gamma between 0.01 to 0.2)

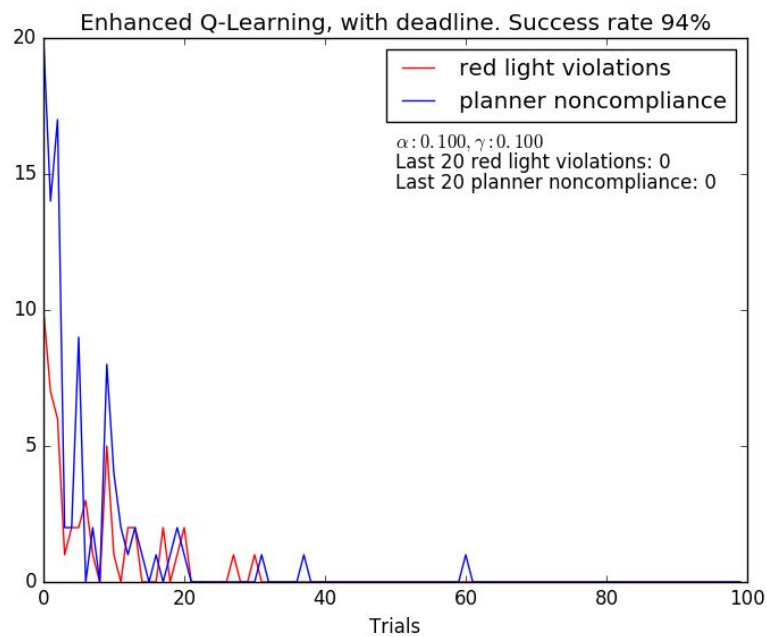
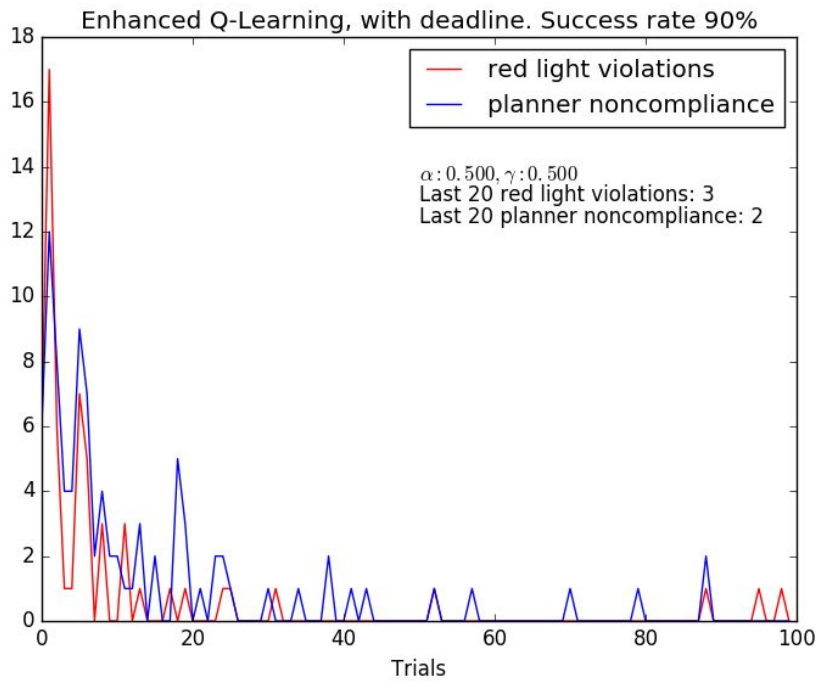


$\alpha$  (alpha), and  $\gamma$  (gamma) - A value of 0.1 chosen as it seemed balanced

$\epsilon$  (epsilon) - The epsilon is initially assigned 1 (maximum exploration) and then reduced at each trial by a factor 0.9. This alone helped get over the 90% success rate.

**Initial Q value** - Different values between 0 and 10 were tried. Having the initial value of 0 give better results consistently.

Below plot shows the improvement in Q-learning after fine tuning the parameters.



In the above sample plot there is no red light violations or planner noncompliance after around 60 trials. The success rate is 94%.

**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?

To confirm the Q-Learning improvements the tests done for the baseline were repeated with the new alpha (0.1) and gamma(0.1) values. The results are as below:

10 tests each with 100 trials

Test #	Success Rate	Red Light Violations in last 20 trials	Planner noncompliances in last 20 trials
1	92	0	1
2	91	0	0
3	95	1	1
4	95	1	0
5	92	4	2
6	96	0	3
7	92	0	1
8	94	0	0
9	94	2	3
10	94	0	0

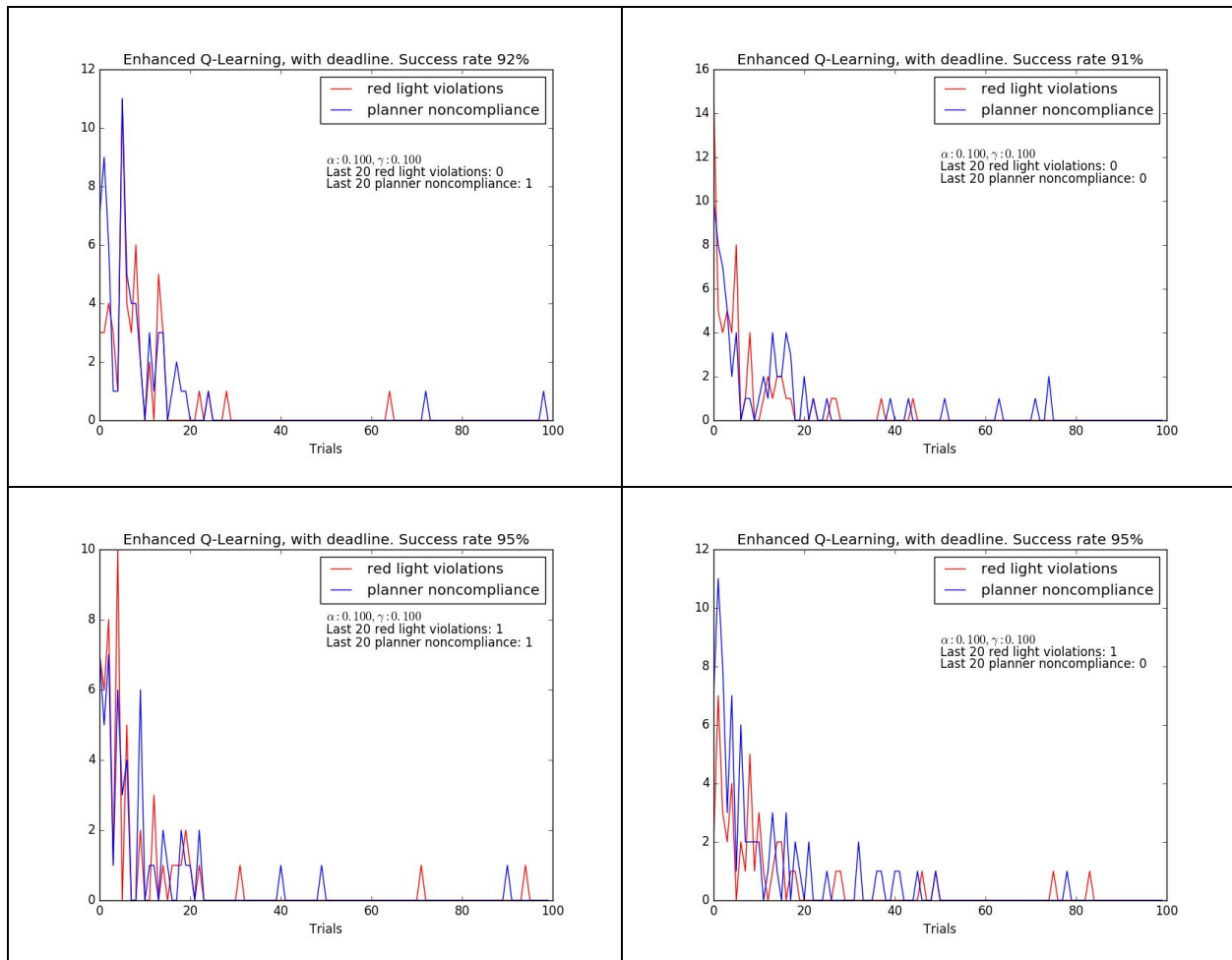
**Summary of Improved Q-Learning(Baseline values in brackets):**

	Success Rate	Red Light Violations in last 20 trials	Planner noncompliances in last 20 trials
<b>Mean</b>	93.5 (88.8)	0.8 (2.7)	1.1 (45.5)
<b>Minimum</b>	91 (69)	0 (0)	0 (2)
<b>Max</b>	96 (98)	4 (7)	3 (114)

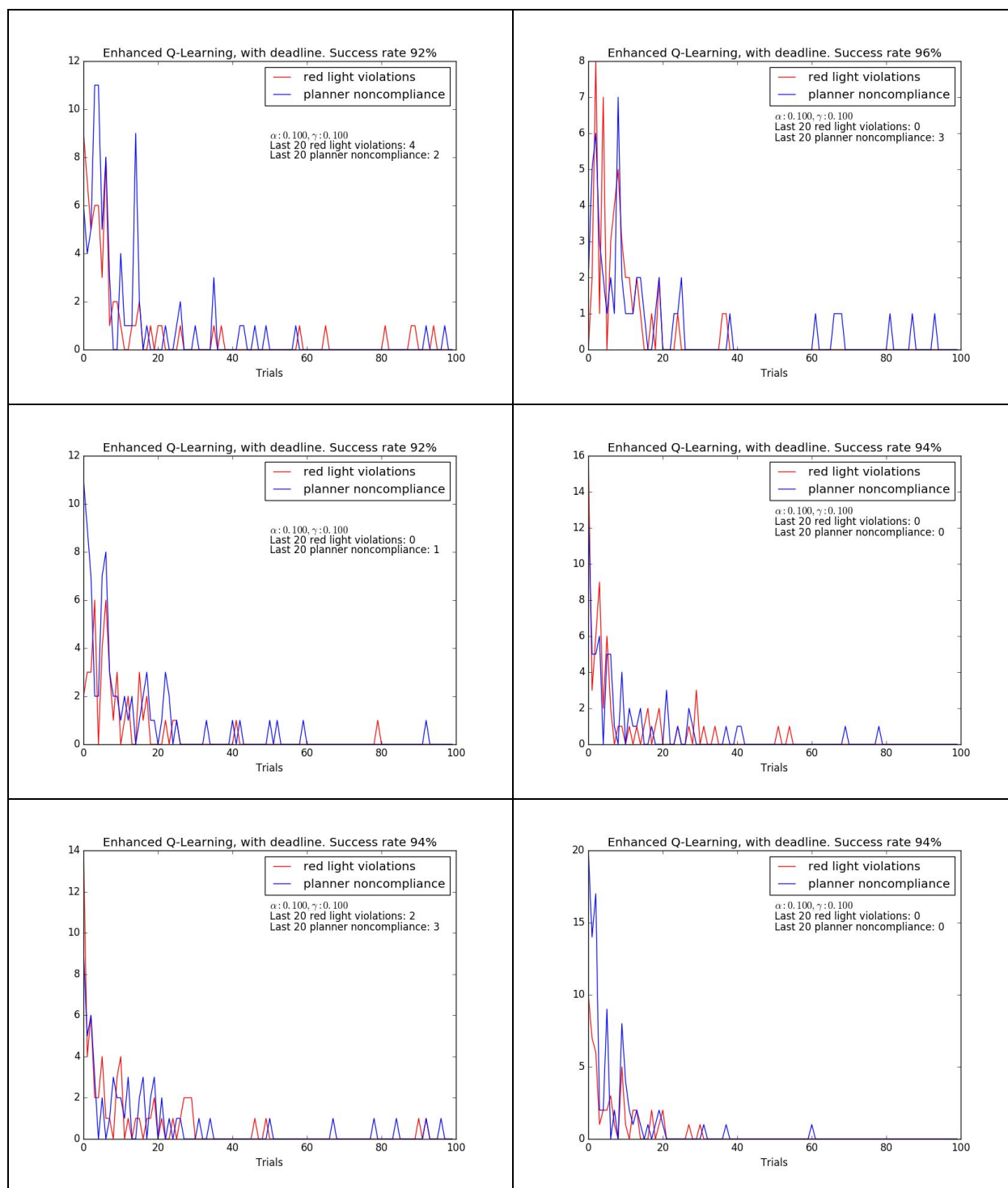
**Observations:**

The overall success rate improved from baseline 88.8 to 93.5. With the improved Q-Learning agent performs more consistently in the range of 91-96% success. For the baseline the range was 69-98%. Significant reduction in red light violations, 2.7 reduced to 0.8. This means the agent has learned to follow the rules of the road. This is evident in the sharp reduction in the red light violations. Tremendous improvement is seen in following the planner suggestions. The average noncompliance to planner got reduced from 45.5 to 1.1. The highest noncompliance recorded with baseline was 114, after enhancement it has come down to 3.

More plots on improved Q-learning below:







Yes, the agent has learned to avoid penalties and to use the optimal route ( route suggested by the planner). A policy simply tells the agent which action to take when they find themselves in the given state. The agent has achieved a success rate of above 93% on average. It learns very fast, there is a sharp reduction in penalties within the first 40 trials itself.

**Sources:**

<https://en.wikipedia.org/wiki/Q-learning>

Bennett, J. Machine learning, part 3: The q-learning algorithm (2016).

<http://articles.wearepop.com/secret-formula-for-self-learning-computers>

The down arrow icon:

<http://www.clipartkid.com/images/42/green-arrow-down-clip-art-at-clker-com-vector-clip-art-online-GcH2ey-clipart.png>