

Machine Learning Engineer Nanodegree

Reinforcement Learning

Project: Train a Smartcab to Drive

Welcome to the fourth project of the Machine Learning Engineer Nanodegree! In this notebook, template code has already been provided for you to aid in your analysis of the *Smartcab* and your implemented learning algorithm. You will not need to modify the included code beyond what is requested. There will be questions that you must answer which relate to the project and the visualizations provided in the notebook. Each section where you will answer a question is preceded by a '**Question X**' header. Carefully read each question and provide thorough answers in the following text boxes that begin with '**Answer:**'. Your project submission will be evaluated based on your answers to each of the questions and the implementation you provide in `agent.py`.

Note: Code and Markdown cells can be executed using the **Shift + Enter** keyboard shortcut. In addition, Markdown cells can be edited by typically double-clicking the cell to enter edit mode.

Getting Started

In this project, you will work towards constructing an optimized Q-Learning driving agent that will navigate a *Smartcab* through its environment towards a goal. Since the *Smartcab* is expected to drive passengers from one location to another, the driving agent will be evaluated on two very important metrics: **Safety** and **Reliability**. A driving agent that gets the *Smartcab* to its destination while running red lights or narrowly avoiding accidents would be considered **unsafe**. Similarly, a driving agent that frequently fails to reach the destination in time would be considered **unreliable**. Maximizing the driving agent's **safety** and **reliability** would ensure that *Smartcabs* have a permanent place in the transportation industry.

Safety and **Reliability** are measured using a letter-grade system as follows:

Grade	Safety	Reliability
A+	Agent commits no traffic violations, and always chooses the correct action.	Agent reaches the destination in time for 100% of trips.
A	Agent commits few minor traffic violations, such as failing to move on a green light.	Agent reaches the destination on time for at least 90% of trips.
B	Agent commits frequent minor traffic violations, such as failing to move on a green light.	Agent reaches the destination on time for at least 80% of trips.
C	Agent commits at least one major traffic violation, such as driving through a red light.	Agent reaches the destination on time for at least 70% of trips.
D	Agent causes at least one minor accident, such as turning left on green with oncoming traffic.	Agent reaches the destination on time for at least 60% of trips.
F	Agent causes at least one major accident, such as driving through a red light with cross-traffic.	Agent fails to reach the destination on time for at least 60% of trips.

To assist evaluating these important metrics, you will need to load visualization code that will be used later on in the project. Run the code cell below to import this code which is required for your analysis.

```
In [1]: # Import the visualization code
import visuals as vs

# Pretty display for notebooks
%matplotlib inline
```

```
In [30]: import numpy as np
import matplotlib.pyplot as plt
```

Understand the World

Before starting to work on implementing your driving agent, it's necessary to first understand the world (environment) which the *Smartcab* and driving agent work in. One of the major components to building a self-learning agent is understanding the characteristics about the agent, which includes how the agent operates. To begin, simply run the `agent.py` agent code exactly how it is -- no need to make any additions whatsoever. Let the resulting simulation run for some time to see the various working components. Note that in the visual simulation (if enabled), the **white vehicle** is the *Smartcab*.

Question 1

In a few sentences, describe what you observe during the simulation when running the default `agent.py` agent code. Some things you could consider:

- *Does the Smartcab move at all during the simulation?*
- *What kind of rewards is the driving agent receiving?*
- *How does the light changing color affect the rewards?*

Hint: From the `/smartcab/` top-level directory (where this notebook is located), run the command

```
'python smartcab/agent.py'
```

Answer: The default agent does not take any actions, irrespective of the state of the environment. This results in positive rewards during the states that precisely require the agent to stay stationary - e.g. when the traffic light is red, or when it should turn left on a green light with oncoming traffic. But it results in much bigger negative rewards when the traffic light is green and fails to move.

Understand the Code

In addition to understanding the world, it is also necessary to understand the code itself that governs how the world, simulation, and so on operate. Attempting to create a driving agent would be difficult without having at least explored the "*hidden*" devices that make everything work. In the `/smartcab/` top-level directory, there are two folders: `/logs/` (which will be used later) and `/smartcab/`. Open the `/smartcab/` folder and explore each Python file included, then answer the following question.

Question 2

- In the agent.py Python file, choose three flags that can be set and explain how they change the simulation.
- In the environment.py Python file, what Environment class function is called when an agent performs an action?
- In the simulator.py Python file, what is the difference between the 'render_text()' function and the 'render()' function?
- In the planner.py Python file, will the 'next_waypoint()' function consider the North-South or East-West direction first?

Answer: There are several arguments that can be set in the agent.py file that change the simulation, here are three examples.

- **num_dummies** Setting this argument when initializing the Environment class allows you to control the amount of traffic in the simulator.
- **grid_size** Setting this argument when initializing the Environment class allows you to control the number of blocks (and in turn the number of intersections) in the simulator.
- **learning** Setting this argument to True when creating the agent using env.create_agent() forces the agent to make use of Q learning.

When an agent performs an action, the act() function of the Environment class is called, which checks that the action is legal and returns a reward.

For the Simulator class there is a render_text() function that prints out feedback about the simulator in the command line. The information includes what action was taken by the agent as well as the reward. The render() function updates the GUI display of the game state.

The next_waypoint() function in the RoutePlanner class it first considers the East-West direction that the agent needs to head in before considering the North-South direction.

Implement a Basic Driving Agent

The first step to creating an optimized Q-Learning driving agent is getting the agent to actually take valid actions. In this case, a valid action is one of None, (do nothing) 'Left' (turn left), 'Right' (turn right), or 'Forward' (go forward). For your first implementation, navigate to the 'choose_action()' agent function and make the driving agent randomly choose one of these actions. Note that you have access to several class variables that will help you write this functionality, such as 'self.learning' and 'self.valid_actions'. Once implemented, run the agent file and simulation briefly to confirm that your driving agent is taking a random action each time step.

Basic Agent Simulation Results

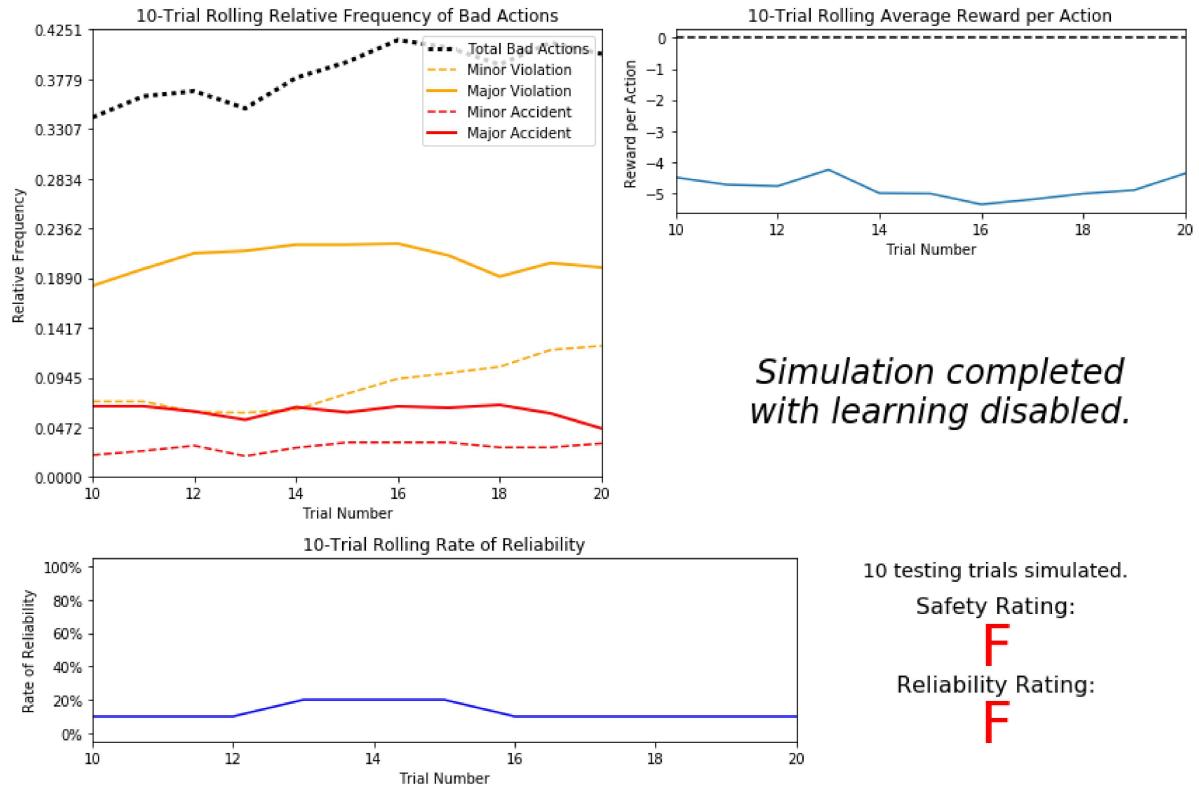
To obtain results from the initial simulation, you will need to adjust following flags:

- 'enforce_deadline' - Set this to True to force the driving agent to capture whether it reaches the destination in time.
- 'update_delay' - Set this to a small value (such as 0.01) to reduce the time between steps in each trial.
- 'log_metrics' - Set this to True to log the simulation results as a .csv file in /logs/.
- 'n_test' - Set this to '10' to perform 10 testing trials.

Optionally, you may disable the visual simulation (which can make the trials go faster) by setting the 'display' flag to False. Flags that have been set here should be returned to their default setting when debugging. It is important that you understand what each flag does and how it affects the simulation!

Once you have successfully completed the initial simulation (there should have been 20 training trials and 10 testing trials), run the code cell below to visualize the results. Note that log files are overwritten when identical simulations are run, so be careful with what log file is being loaded! Run the agent.py file after setting the flags from projects/smartcab folder instead of projects/smartcab/smartcab.

```
In [3]: # Load the 'sim_no-Learning' Log file from the initial simulation results
vs.plot_trials('sim_no-learning.csv')
```



Question 3

Using the visualization above that was produced from your initial simulation, provide an analysis and make several observations about the driving agent. Be sure that you are making at least one observation about each panel present in the visualization. Some things you could consider:

- *How frequently is the driving agent making bad decisions? How many of those bad decisions cause accidents?*
- *Given that the agent is driving randomly, does the rate of reliability make sense?*
- *What kind of rewards is the agent receiving for its actions? Do the rewards suggest it has been penalized heavily?*
- *As the number of trials increases, does the outcome of results change significantly?*
- *Would this Smartcab be considered safe and/or reliable for its passengers? Why or why not?*

Answer: The agent is making bad decisions at a rate of about 33-42% of the time, approximately 20-30% of those mistakes lead to accidents.

The rate of reliability is quite low, ranging from about 10-20%, which is not surprising given that it is randomly choosing actions.

The agent is receiving net negative rewards for each trial, of approximately -4 to -5.

The agent is not improving over time, it does not get better rewards or reliability.

This agent would certainly be considered very unsafe, it causes too many accidents. It is also very unreliable. There is no consumer that would pay any money for a taxi that arrives at its intended destination less than 20% of the time and gets into a serious accident 7% of the time. The taxi company would go out of business.

Inform the Driving Agent

The second step to creating an optimized Q-learning driving agent is defining a set of states that the agent can occupy in the environment. Depending on the input, sensory data, and additional variables available to the driving agent, a set of states can be defined for the agent so that it can eventually *learn* what action it should take when occupying a state. The condition of '*if state then action*' for each state is called a **policy**, and is ultimately what the driving agent is expected to learn. Without defining states, the driving agent would never understand which action is most optimal -- or even what environmental variables and conditions it cares about!

Identify States

Inspecting the 'build_state()' agent function shows that the driving agent is given the following data from the environment:

- 'waypoint', which is the direction the *Smartcab* should drive leading to the destination, relative to the *Smartcab's* heading.
- 'inputs', which is the sensor data from the *Smartcab*. It includes
 - 'light', the color of the light.
 - 'left', the intended direction of travel for a vehicle to the *Smartcab's* left. Returns None if no vehicle is present.
 - 'right', the intended direction of travel for a vehicle to the *Smartcab's* right. Returns None if no vehicle is present.
 - 'oncoming', the intended direction of travel for a vehicle across the intersection from the *Smartcab*. Returns None if no vehicle is present.
- 'deadline', which is the number of actions remaining for the *Smartcab* to reach the destination before running out of time.

Question 4

*Which features available to the agent are most relevant for learning both **safety** and **efficiency**? Why are these features appropriate for modeling the Smartcab in the environment? If you did not choose some features, why are those features not appropriate?*

Answer: The inputs is certainly the most important information for the safety of the smart-cab. Without this information it has no way of knowing the traffic conditions and will lead to accidents. However it does not need all of the inputs data. It can do without the right information due to the USA right of way traffic rules, in which cars on the right are not relevant for making any decisions. The waypoint information would be the most relevant information for making sure the agent actually goes in the correct direction to reach its destination. Deadline would be the least relevant information for making decisions as to which direction to move in, and would also increase the state space unnecesarially.

Define a State Space

When defining a set of states that the agent can occupy, it is necessary to consider the size of the state space. That is to say, if you expect the driving agent to learn a **policy** for each state, you would need to have an optimal action for every state the agent can occupy. If the number of all possible states is very large, it might be the case that the driving agent never learns what to do in some states, which can lead to uninformed decisions. For example, consider a case where the following features are used to define the state of the *Smartcab*:

```
('is_raining', 'is_foggy', 'is_red_light', 'turn_left', 'no_traffic',
'previous_turn_left', 'time_of_day').
```

How frequently would the agent occupy a state like (False, True, True, True, False, False, '3AM')? Without a near-infinite amount of time for training, it's doubtful the agent would ever learn the proper action!

Question 5

*If a state is defined using the features you've selected from **Question 4**, what would be the size of the state space? Given what you know about the environment and how it is simulated, do you think the driving agent could learn a policy for each possible state within a reasonable number of training trials?*

Hint: Consider the combinations of features to calculate the total number of states!

Answer: The four features that were selected have the following possible values.

feature	possible values
'light'	'red', 'green'
'left'	None, 'forward', 'left', 'right'
'oncoming'	None, 'forward', 'left', 'right'
'waypoint'	None, 'forward', 'left', 'right'

Even though there are 4 possible states for waypoint, the only time that None will appear as a value is when it reaches its destination. We can therefore treat it as a feature with three possible values when calculating the statespace:

$$2 \times 4 \times 4 \times 3 = 96$$

Theoretically the minimum number of steps that could potentially cover all possible states is 96. But to explore all possible *actions* for all the possible states would require a theoretical minimum of 384 steps. Of course in practice there will be many states that are much more common, and some that are much more rare, so it would require running for much more than 384 time steps.

Assuming somewhere in the vicinity of 20 time steps for each trial, it would require a theoretical minimum of about 20 trials to cover all states and actions. In practice, therefore it will require several hundred or even thousands of trials to explore all states and actions.

Update the Driving Agent State

For your second implementation, navigate to the 'build_state()' agent function. With the justification you've provided in **Question 4**, you will now set the 'state' variable to a tuple of all the features necessary for Q-Learning. Confirm your driving agent is updating its state by running the agent file and simulation briefly and note whether the state is displaying. If the visual simulation is used, confirm that the updated state corresponds with what is seen in the simulation.

Note: Remember to reset simulation flags to their default setting when making this observation!

Implement a Q-Learning Driving Agent

The third step to creating an optimized Q-Learning agent is to begin implementing the functionality of Q-Learning itself. The concept of Q-Learning is fairly straightforward: For every state the agent visits, create an entry in the Q-table for all state-action pairs available. Then, when the agent encounters a state and performs an action, update the Q-value associated with that state-action pair based on the reward received and the iterative update rule implemented. Of course, additional benefits come from Q-Learning, such that we can have the agent choose the *best* action for each state based on the Q-values of each state-action pair possible. For this project, you will be implementing a *decaying*, ϵ -greedy Q-learning algorithm with *no* discount factor. Follow the implementation instructions under each **TODO** in the agent functions.

Note that the agent attribute `self.Q` is a dictionary: This is how the Q-table will be formed. Each state will be a key of the `self.Q` dictionary, and each value will then be another dictionary that holds the *action* and *Q-value*. Here is an example:

```
{
    'state-1': {
        'action-1' : Qvalue-1,
        'action-2' : Qvalue-2,
        ...
    },
    'state-2': {
        'action-1' : Qvalue-1,
        ...
    },
    ...
}
```

Furthermore, note that you are expected to use a *decaying ϵ (exploration) factor*. Hence, as the number of trials increases, ϵ should decrease towards 0. This is because the agent is expected to learn from its behavior and begin acting on its learned behavior. Additionally, The agent will be tested on what it has learned after ϵ has passed a certain threshold (the default threshold is 0.01). For the initial Q-Learning implementation, you will be implementing a linear decaying function for ϵ .

Q-Learning Simulation Results

To obtain results from the initial Q-Learning implementation, you will need to adjust the following flags and setup:

- 'enforce_deadline' - Set this to True to force the driving agent to capture whether it reaches the destination in time.
- 'update_delay' - Set this to a small value (such as 0.01) to reduce the time between steps in each trial.
- 'log_metrics' - Set this to True to log the simulation results as a .csv file and the Q-table as a .txt file in /logs/.
- 'n_test' - Set this to '10' to perform 10 testing trials.
- 'learning' - Set this to 'True' to tell the driving agent to use your Q-Learning implementation.

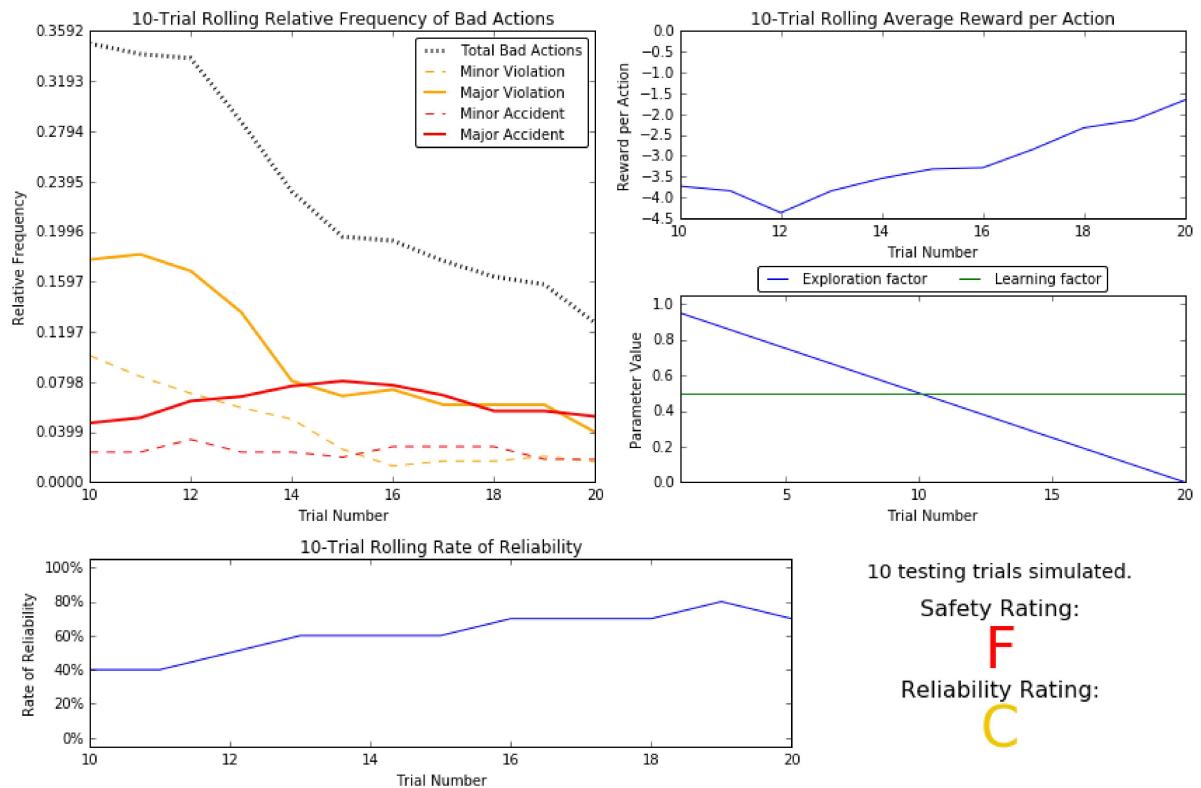
In addition, use the following decay function for ϵ :

$$\epsilon_{t+1} = \epsilon_t - 0.05, \text{ for trial number } t$$

If you have difficulty getting your implementation to work, try setting the 'verbose' flag to True to help debug. Flags that have been set here should be returned to their default setting when debugging. It is important that you understand what each flag does and how it affects the simulation!

Once you have successfully completed the initial Q-Learning simulation, run the code cell below to visualize the results. Note that log files are overwritten when identical simulations are run, so be careful with what log file is being loaded!

```
In [2]: # Load the 'sim_default-learning' file from the default Q-Learning simulation
vs.plot_trials('sim_default-learning.csv')
```



Question 6

Using the visualization above that was produced from your default Q-Learning simulation, provide an analysis and make observations about the driving agent like in **Question 3**. Note that the simulation should have also produced the Q-table in a text file which can help you make observations about the agent's learning. Some additional things you could consider:

- *Are there any observations that are similar between the basic driving agent and the default Q-Learning agent?*
- *Approximately how many training trials did the driving agent require before testing? Does that number make sense given the epsilon-tolerance?*
- *Is the decaying function you implemented for ϵ (the exploration factor) accurately represented in the parameters panel?*
- *As the number of training trials increased, did the number of bad actions decrease? Did the average reward increase?*
- *How does the safety and reliability rating compare to the initial driving agent?*

Answer: The agent trains for about 20 trials which is what would be expected given a linear decay of 0.05 and an epsilon threshold of 0.05. We can see this linear decay quite clearly in the plot. We can see that as the number of trials increased, the relative frequency of traffic violations went down, however, the relative number of accidents stays fairly constant. The rewards per action and the reliability also went up. The final reliability rating improved, compared to the initial driving agent, with a score of C. The safety rating however is still an F.

Improve the Q-Learning Driving Agent

The third step to creating an optimized Q-Learning agent is to perform the optimization! Now that the Q-Learning algorithm is implemented and the driving agent is successfully learning, it's necessary to tune settings and adjust learning parameters so the driving agent learns both **safety** and **efficiency**. Typically this step will require a lot of trial and error, as some settings will invariably make the learning worse. One thing to keep in mind is the act of learning itself and the time that this takes: In theory, we could allow the agent to learn for an incredibly long amount of time; however, another goal of Q-Learning is to *transition from experimenting with unlearned behavior to acting on learned behavior*. For example, always allowing the agent to perform a random action during training (if $\epsilon = 1$ and never decays) will certainly make it *learn*, but never let it *act*. When improving on your Q-Learning implementation, consider the implications it creates and whether it is logically sensible to make a particular adjustment.

Improved Q-Learning Simulation Results

To obtain results from the initial Q-Learning implementation, you will need to adjust the following flags and setup:

- 'enforce_deadline' - Set this to True to force the driving agent to capture whether it reaches the destination in time.
- 'update_delay' - Set this to a small value (such as 0.01) to reduce the time between steps in each trial.
- 'log_metrics' - Set this to True to log the simulation results as a .csv file and the Q-table as a .txt file in /logs/.
- 'learning' - Set this to 'True' to tell the driving agent to use your Q-Learning implementation.
- 'optimized' - Set this to 'True' to tell the driving agent you are performing an optimized version of the Q-Learning implementation.

Additional flags that can be adjusted as part of optimizing the Q-Learning agent:

- 'n_test' - Set this to some positive number (previously 10) to perform that many testing trials.
- 'alpha' - Set this to a real number between 0 - 1 to adjust the learning rate of the Q-Learning algorithm.
- 'epsilon' - Set this to a real number between 0 - 1 to adjust the starting exploration factor of the Q-Learning algorithm.
- 'tolerance' - set this to some small value larger than 0 (default was 0.05) to set the epsilon threshold for testing.

Furthermore, use a decaying function of your choice for ϵ (the exploration factor). Note that whichever function you use, it **must decay to 'tolerance' at a reasonable rate**. The Q-Learning agent will not begin testing until this occurs. Some example decaying functions (for t , the number of trials):

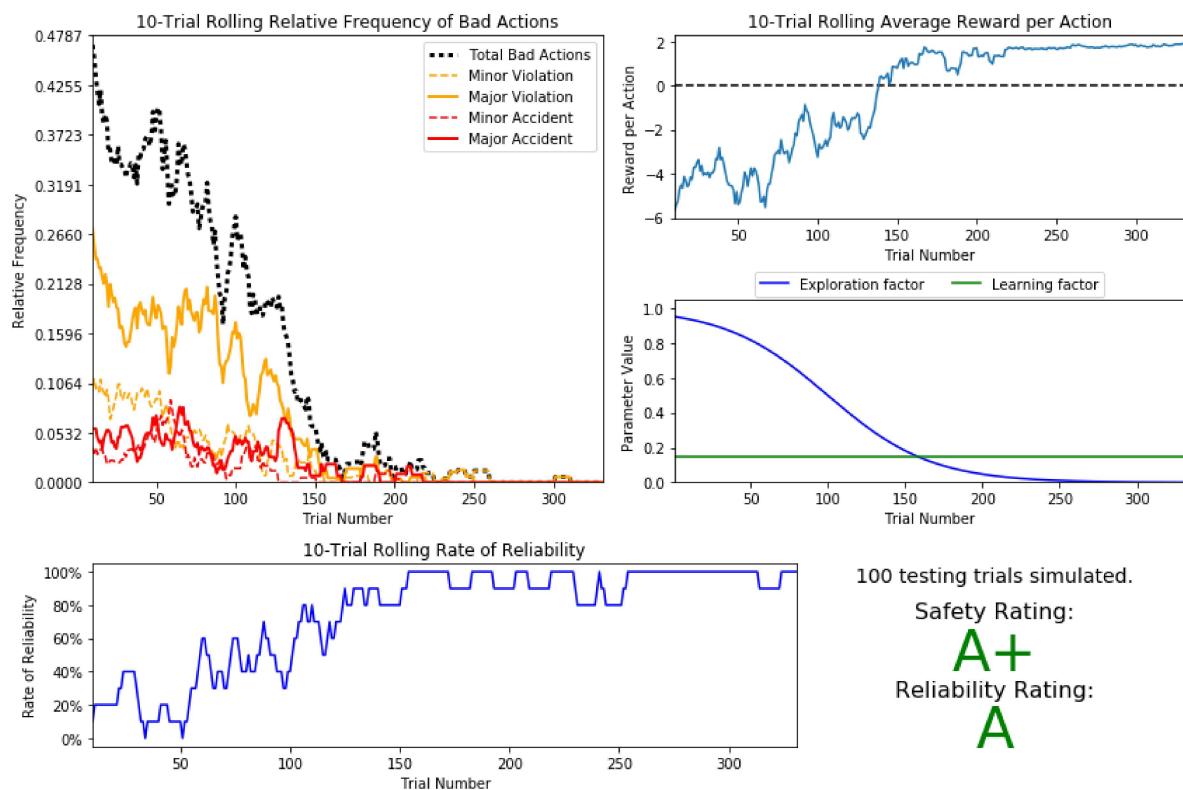
$$\epsilon = a^t, \text{ for } 0 < a < 1 \quad \epsilon = \frac{1}{t^2} \quad \epsilon = e^{-at}, \text{ for } 0 < a < 1 \quad \epsilon = \cos(at), \text{ for } 0 < a < 1$$

You may also use a decaying function for α (the learning rate) if you so choose, however this is typically less common. If you do so, be sure that it adheres to the inequality $0 \leq \alpha \leq 1$.

If you have difficulty getting your implementation to work, try setting the 'verbose' flag to True to help debug. Flags that have been set here should be returned to their default setting when debugging. It is important that you understand what each flag does and how it affects the simulation!

Once you have successfully completed the improved Q-Learning simulation, run the code cell below to visualize the results. Note that log files are overwritten when identical simulations are run, so be careful with what log file is being loaded!

```
In [29]: # Load the 'sim_improved-Learning' file from the improved Q-Learning simulation
vs.plot_trials('sim_improved-learning.csv')
```



Question 7

Using the visualization above that was produced from your improved Q-Learning simulation, provide a final analysis and make observations about the improved driving agent like in **Question 6**. Questions you should answer:

- What decaying function was used for epsilon (the exploration factor)?
- Approximately how many training trials were needed for your agent before begining testing?
- What epsilon-tolerance and alpha (learning rate) did you use? Why did you use them?
- How much improvement was made with this Q-Learner when compared to the default Q-Learner from the previous section?
- Would you say that the Q-Learner results show that your driving agent successfully learned an appropriate policy?
- Are you satisfied with the safety and reliability ratings of the Smartcab?

Answer:

Epsilon was decayed using the following function:

$$1 - \frac{1}{1 + e^{-k(t-offset)}}$$

Which is an inverse sigmoid function, such that:

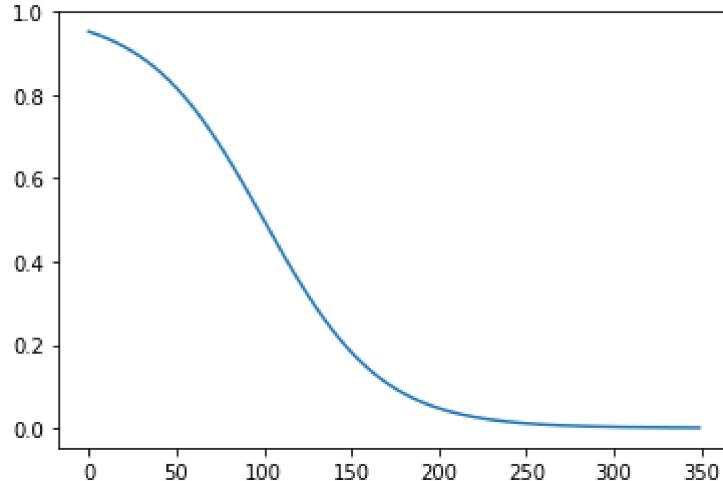
- **t**: is the trial number
- **k**: controls the rate of decay (this was set to 0.03)
- **offset**: shifts the sigmoid, controlling at what point along the curve the value will be at time step 0 (this was set to 100).

This decay function can be visualised in the following plot.

```
In [39]: k = 0.03
offset = 100

x = np.arange(0,350)
y = 1 - (1. / (1. + np.exp(-k * (x-offset))))
plt.plot(x,y)
```

```
Out[39]: []
```



The rationale behind this decay rate for epsilon decay was that it would force the model to explore a lot more actions during the beginning, so that it would populate the Q table with rules of thumb, and rapidly decay so that it could start exploiting the learnt values, and rely less on exploration.

Alpha was kept constant at a rate of 0.15

The epsilon tolerance was reduced to 0.001 to allow it to continue training for longer.

There was a significant improvement with the setting used for this Q-Learner compared to the default Q-learner. On the test trials, the agent gets to its destination at least 90% of the time, without making any traffic violations or getting into accidents. This indicates that the agent has learnt a very good policy, at least for the states that existed in the testing trials, and means it is quite reliable and safe under those states.

Define an Optimal Policy

Sometimes, the answer to the important question "*what am I trying to get my agent to learn?*" only has a theoretical answer and cannot be concretely described. Here, however, you can concretely define what it is the agent is trying to learn, and that is the U.S. right-of-way traffic laws. Since these laws are known information, you can further define, for each state the *Smartcab* is occupying, the optimal action for the driving agent based on these laws. In that case, we call the set of optimal state-action pairs an **optimal policy**. Hence, unlike some theoretical answers, it is clear whether the agent is acting "incorrectly" not only by the reward (penalty) it receives, but also by pure observation. If the agent drives through a red light, we both see it receive a negative reward but also know that it is not the correct behavior. This can be used to your advantage for verifying whether the **policy** your driving agent has learned is the correct one, or if it is a **suboptimal policy**.

Question 8

Provide a few examples (using the states you've defined) of what an optimal policy for this problem would look like. Afterwards, investigate the 'sim_improved-learning.txt' text file to see the results of your improved Q-Learning algorithm. *For each state that has been recorded from the simulation, is the **policy** (the action with the highest value) correct for the given state? Are there any states where the policy is different than what would be expected from an optimal policy?* Provide an example of a state and all state-action rewards recorded, and explain why it is the correct policy.

Answer:

For the following state, where there is a green light, and the destination is to the left, and there is no other traffic, then the optimal policy should be to turn **left**.

light	oncomming	traffic from left	waypoint
'green'	None	None	'left'

The Qvalues accurately capture this, it even captures that the most incorrect thing to do is to stay still, which would be a violation of the traffic rules.

```
('green', None, None, 'left')
-- forward : 0.80
-- left : 2.00
-- right : 0.97
-- None : -5.08
```

For the following state, the optimal action should be to turn **right**, because according to US traffic rules you are allowed to turn right, even if you are at a red traffic light, as long as there is no traffic coming along the lane you want to turn into. The agent has learned that the most correct thing to do is to turn right.

light	oncomming	traffic from left	waypoint
'red'	None	None	'right'

The Q values for that state reflect this optimal policy quite nicely.

```
('red', None, None, 'right')
-- forward : -19.65
-- left : -11.51
-- right : 2.22
-- None : 1.52
```

For the following state where there is a green traffic light, and the destination is to the left, but there is oncomming traffic that is moving forward, then the optimal policy should be for the car to **wait** for oncomming traffic to clear out before turning left.

light	oncomming	traffic from left	waypoint
'green'	'forward'	None	'left'

The policy learnt by the agent is an example of a suboptimal policy. It does capture that the most incorrect thing to do would be to turn left since it would not only violate the traffic rules, but also cause an accident. However, it is also penalizing the optimal policy quite heavily, and the maximum Q value is assigned to the action of going forward instead, which might still get the car to its intended destination, but it is not the optimal solution.

```
('green', 'forward', None, 'left')
-- forward : 0.71
-- left : -16.79
-- right : 0.24
-- None : -4.40
```