

COMP5212 Machine Learning 2018 Fall programming project proposal

Building a self-driving agent with Reinforcement Learning

Hok Chun Ng, 20272532, hcngac@connect.ust.hk
Shengyuan Zhang, 20565161, szhangcg@connect.ust.hk
Ge Chen, 20360858, gchenaj@connect.ust.hk

October 16, 2018

Abstract

In this project proposal, we intend to build a self-driving agent that only relies on vision using reinforcement learning algorithms. Due to the limitation of real-world hardware and data, we intend to use a game simulator provided by OpenAI [3] to generate training data. We introduce the elements of Q-learning, a reinforcement learning technique which does not require a model of the environment.

1 Application and practical significance

Self-driving car has been a hot topic in recent years. Several industry giants, like Google [1] and Tesla Motor [2], have devoted significant efforts to the developing of self-driving cars.

Applications of self-driving agent can be wide. Example includes long-distance truck driving, smart taxi and bus, or even be incorporated into large-scale autonomous traffic systems.

Letting computers to drive can release human beings from the boring and tiring job of driving as well as reduce the frequency of traffic accidents. It also opens up new business and city-management opportunities including smart traffic system and smart taxi services. Applications are wide and undiscovered and the significance of perfecting self-driving agent is huge in opening the door to these huge possibilities.

However, designing a robust self-driving system is non-trivial as the real-world traffic conditions are diversified. Limited by hardware aquirement and law, we intend to find a solution to self-driving in a simulated game environment in the hope that the result can be transfered to real-world self-driving.

2 Problem formulation

We formulate the problem of self-driving as a continuous decision making process. The self-driving agent receives input from sensors, evaluate each possible actions and choose the best action.

For human drivers, the driving behavior can be formulated in a loop from an environment sensing input to action output. Our eyes and ears are the sensors to interpret the current environment, including the road view from wind screen (e.g. weather, traffic lights, walking people), view from side mirror (e.g. neighboring cars, following cars), car horns, and so on. Then our brain takes all these input signals to decide what action to take. The actions include turning the steering wheel, lighting turn signals, accelerating, braking, and so on. Once the actions are executed, the environment input changes, and again drivers will decide and take a new round of actions accordingly. This loop continues until the car arrives at the destination.

We try to mimic how human decision works with machine learning algorithms. The machine learning problem in self-driving can be: given a set of environment input, find the best actions to execute until arriving at destination.

3 Data set

In this project, we intend to use the OpenAI Gym library to simulate the environment of car driving. The OpenAI Gym [3] is a toolkit for developing and comparing reinforcement algorithms. It is a collection of test problems, by providing the environments that anyone can use to test her own reinforcement learning algorithms. The environments have a shared interface, allowing any user to write general algorithms. It is a good tool as a common ground for benchmarking of machine learning models.

We will focus on the environment “Enduro-v0” [4], which is the Atari 2600 game Enduro simulating a driving experience. In this environment, the observation is an RGB image of the screen (simulate the view from wind screen), which is an array of shape (210, 160, 3). Each action is repeatedly performed for a duration of k frames, where k is uniformly sampled from $\{2, 3, 4\}$. To reduce computation overhead, we may first converse the RGB images to grayscale ones. The reward is the score of the game, which is related to the distance the car goes.

4 Machine learning methods

The environment-action loop in the driving experience resembles the basics of reinforcement learning, an area of machine learning concerned with taking actions in an environment in order to maximize some notion of cumulative reward. Unlike supervised learning, reinforcement learning does not require correct input/output pairs. Instead, reinforcement learning often requires a reward function in terms of different actions under the environment. For driving, the final reward can be no traffic accident and arrive at the destination safely.

Figure 1 shows the basic model of reinforcement learning [5]. At each time t , an agent receives the environment state S_t and current reward R_t . It then

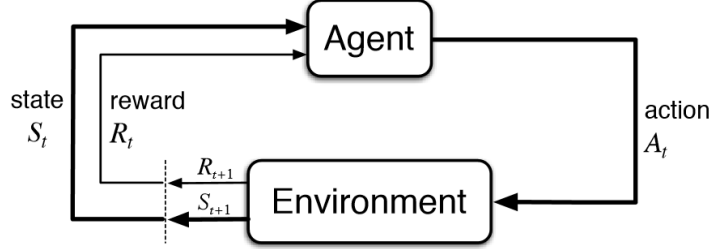


Figure 1: Reinforcement Learning Illustration [5]

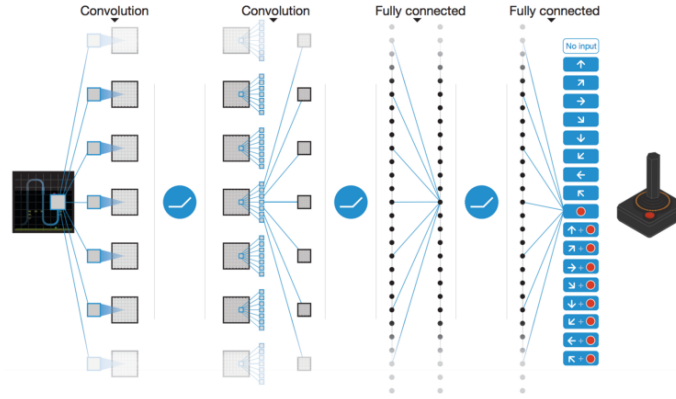


Figure 2: Deep Q-learning Network Illustration [6]

chooses an action a_t from the set of available actions, which is sent to the environment. The environment moves to a new state S_{t+1} and new reward R_{t+1} and feedback to the agent. The goal of a reinforcement learning agent is to gain as much as reward as possible.

In this proposed project, we will use the method of deep Q-learning. Q-learning defines a reward expectation function, Q , which provides the expected reward of an action a taken under the current state s . State is the image displayed and actions are all available actions of the drive. This Q function can be recursively defined as $Q(s, a) = R(s, a) + \alpha \max_{a'} Q(S'(s, a), a')$, where R is the immediate reward given action a taken under state s , α is the discount rate of future reward and S' is the state transition. Given numerous possibilities of imagery state s , we cannot compute Q directly. Alternatively, we can use a deep convolutional network to estimate Q .

Figure 2 illustrates a model of deep Q-learning network [6]. First, we have several layers of convolution and max-pooling to reduce the imagery input into smaller vector. Then we have several fully-connected layer for the final decision making. The sample network shown estimates a function $Q : S \rightarrow A^n$, which takes a image input and output the expected reward for all possible actions.

5 Experiment and performance evaluation

Our experiment involves deep Q network training and performance evaluation.

Training is a process of random discovery. The agent is allowed to randomly choose actions in the game until game over. States, action and reward is recorded in the process and is used to train the deep Q network. We evaluate the performance in fixed game intervals to show the learning progress. Training stops when the performance does not increase over certain extend.

In performance evaluation, we let the self-driving agent control the car according to the output of the deep Q network. The agent feed the image display to the network and choose the action with the highest output value in each step, until game over. We run the process for a number of times to get an average performance in each evaluation. Performance is defined as the score obtained at game over.

6 Project planning

Timeline	Project Task
Oct 25	system setup
Oct 30	data generation and preprocessing
Nov 15	algorithm
Nov 25	evaluation
Nov 30	report and video

References

- [1] Google Self-Driving Car, <https://www.google.com/selfdrivingcar>
- [2] Tesla Self-Driving Hardware, <https://www.tesla.com/autopilot>
- [3] OpenAI Gym toolkit, <https://gym.openai.com>
- [4] OpenAI Gym Enduro Environment, <https://gym.openai.com/envs/Enduro-v0>
- [5] Sutton, Richard S., Andrew G. Barto, and Francis Bach. **Reinforcement learning: An introduction**. MIT press, 1998.
- [6] Volodymyr Mnih, Koray Kavukcuoglu, et al. **Human-level control through deep reinforcement learning**. Nature, 2015.