

Exercise 2: A Reactive Agent for the Pickup and Delivery Problem

Group №70: Julien Perrenoud, Pierre-Antoine Desplaces

October 9, 2018

1 Problem Representation

1.1 Representation Description

States - We defined our states by answering the question "what variable changes at each decision point?" with two things: where the agent is and what he knows about the package. The later only includes the destination of the package as its value is already present in the rewards table later. Therefore,

$$\text{States} = \underbrace{\{ \text{City}_1, \dots, \text{City}_n \}}_{\text{Location}} \times \left(\underbrace{\{ \text{NoPackage} \} \cup \{ \text{City}_1, \dots, \text{City}_n \}}_{\text{Delivery}} \right)$$

This representation does not explicitly exclude states where the destination and origin city of a package is the same. Instead, we expect a transition probability of 0 from the system if such states are not permitted.

Actions - This time, we tried to answer the question "what choice does the agent have asked to make a decision?". This is simpler, as our agent can only deliver the package or move to a neighboring city.

$$\text{Actions} = \{ \text{Pickup} \} \cup \{ \text{Move}(\text{City}_1), \dots, \text{Move}(\text{City}_n) \}$$

In the problem description, only a subset of actions are available from each state (as the agent can only move to *neighboring* cities, and an agent cannot pick-up a non-existing task). However, it seemed more natural for us to define the action set independently from the current state in which the agent is in and rely on other mechanisms (such as the rewards table) to explicitly remove impossible actions for each state.

Transition Probability - The first thing to realize when filling the transition table is that a lot of triplets (s, a, s') are actually inconsistent with the model (e.g. if the "Location" attribute of s' is different than the "Delivery" (or "Move") destination of s (or a)). Such cases will always return a probability of 0. For consistent cases, the transition probability depends on the probability of getting a task (or not) from $s'.\text{Location}$ to $s'.\text{Delivery}$ (here defined as p , which also returns the probability of not receiving a package in s' if $s'.\text{Delivery} = \text{NoPackage}$).

$$Pr[s' \mid s, a] = \begin{cases} p(s'.\text{Location}, s'.\text{Delivery}) & \text{if } a = \text{Pickup}, s.\text{Delivery} = s'.\text{Location} \\ p(s'.\text{Location}, s'.\text{Delivery}) & \text{if } a = \text{Move}(\text{City}_k), \text{City}_k = s'.\text{Location} \\ 0, & \text{otherwise} \end{cases}$$

Rewards Table - As we try to optimize the profit per action of our agent, we need to fill our rewards table with the average profit of each action at each particular step. In order to calculate this, we simply use the expected reward of a particular task (if the action is "Pickup" and there is a task, otherwise we use 0) minus the cost associated with the selected travel.

Additionally, we use the rewards table to explicitly mark impossible actions we defined earlier. In order to do this and completely discourage our agent of undertaking such actions, we simply associated them with a reward of $-\infty$. Therefore,

$$R(s, a) = \begin{cases} r(s.\text{Location}, s.\text{Delivery}) - c(s.\text{Location}, s.\text{Delivery}) & \text{if } a = \text{Pickup}, s.\text{Delivery} \neq \text{NoPackage} \\ -c(s.\text{Location}, \text{City}_k) & \text{if } a = \text{Move}(\text{City}_k), \text{City}_k \in N(s.\text{Location}) \\ -\infty & \text{otherwise} \end{cases}$$

where $r(\text{City}_i, \text{City}_j)$ is the *expected* reward of a task from City_i to City_j , $c(\text{City}_i, \text{City}_j)$ is the distance (in kilometers) between City_i and City_j multiplied by the cost per kilometer of the vehicle, and $N(\text{City}_i)$ is the set of neighbors of City_i .

1.2 Implementation Details

States and Actions - In our implementation, we first created `State` and `AgentAction` classes. In order to avoid redundant properties, we set `State.destinationCity` to `null` to indicate that there is no package (this is also consistent with Logist’s `probability` function). Similarly, `AgentAction`’s `moveCity` is `null` to indicate a "Pickup" action.

Rewards, Probabilities and Costs - In our implementation, we replace $p(i, j)$ with `td.probability(i, j)`, $r(i, j)$ with `td.reward(i, j)` and $c(i, j)$, with `i.distanceTo(j) * vehicle.costPerKm()`.

Skip Impossible Actions - In the basic version of `iterateQ` (our version of the MDP Learning algorithm), we iterated over a lot of impossible actions. In order to prevent this and optimize the learning, we simply check and skip the iteration if the reward for the current state and action is `NEGATIVE_INFINITY`. This allows us to skip approximately two thirds of iterations with no change in the final Q-Table.

2 Results

2.1 Experiment 1: Discount factor

2.1.1 Setting

In this first experiment, we test 4 different discount factors [0.25, 0.5, 0.75, 0.95] over 2 simulations (France and Switzerland) for 3000 time steps. The reward per km for each agent can be seen on Figure 1 and their final average reward per action on Table 1.

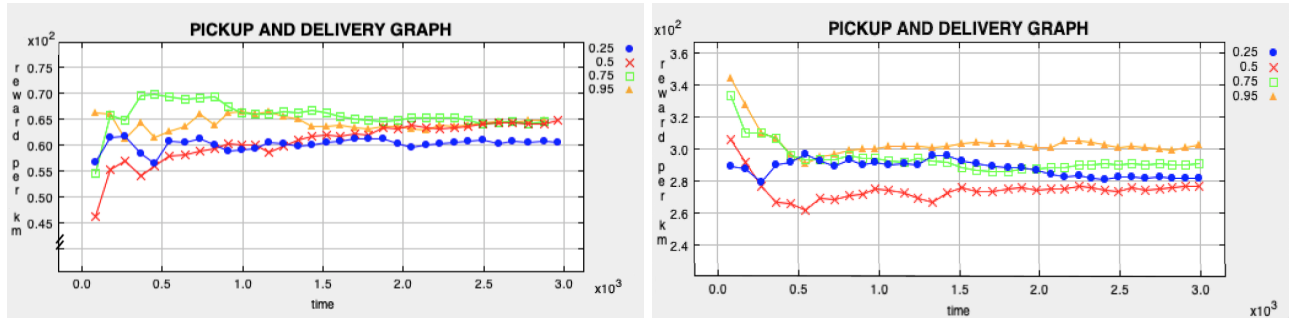


Figure 1: Reward per km over time for agents with different discount factors. On the left, in France and on the right, in Switzerland.

2.1.2 Observations

For the first metric, reward per km, all agents of the first simulation get very similar results, except the one with a factor of 0.25 which is worse, and in the second simulation, the agent with the highest value (0.95) got the best performance. As for the average reward per action (which is the metric we actually try to optimize), the agents are also quite close, but the highest value of 0.95 obtains the best performance in both simulations. For this reason, we decided to select 0.95 as our winner.

Discount \ Country	France	Switzerland
0.25	36339.81	47303.50
0.5	38854.97	46661.04
0.75	38705.52	47087.54
0.95	39031.59	48065.10

Table 1: Experiment 1 - Average reward per action after 3000 time steps for agents with different discount factors.

Agent \ Country	France	Switzerland
Random	31093.81	37188.36
Greedy	36944.03	43974.26
RLA	39138.37	47759.41

Table 2: Experiment 2 - Average reward per action after 3000 time steps for agents with different strategies.

2.2 Experiment 2: Comparisons with dummy agents

2.2.1 Setting

In the second experiment, we compare our reinforcement learning algorithm against 2 different "dummy" behaviors. This is done on the same maps, again for 3000 time steps. The reward per km for each agent can be seen on Figure 2 and their final average reward per action on Table 2. Finally, the agent behaviors are the following:

Random : Accept or deny tasks with a 0.5 probability.

Greedy : Accept every task proposed and otherwise move to the neighboring city maximizing the expected value $\mathbb{E}_{city} = p_{city}(task) * \mathbb{E}_R(city)$.

Reinforcement Learning : Our previous winner, the RLA behavior with a discount factor of 0.95.

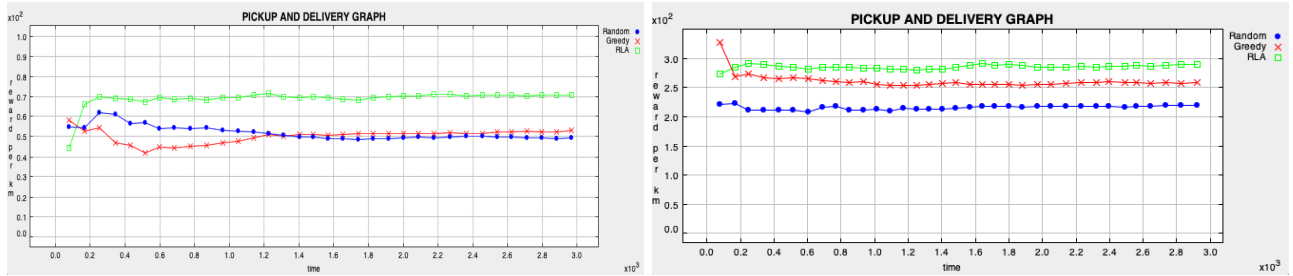


Figure 2: Reward per km over time for agents with different strategies. On the left, in France and on the right, in Switzerland.

2.2.2 Observations

In this experiment, we can see that our agent (RLA) outperforms the greedy and random agents by a clear margin on both metrics, and in both simulations.