

Autonomous Software Agents project

Bozzo Francesco, Izzo Federico University of Trento
Trento, Italy
francesco.bozzo@studenti.unitn.it, federico.izzo@studenti.unitn.it

Index Terms—Autonomous Software Agents; Single-agent; Multi-agent; BDI; A Star; Heuristic

I. AGENT

In the field of *Computer Science* an *agent* is an individual situated in some environment, and capable of flexible autonomous action in that environment in order to meet its design objectives. With respect to a more traditional algorithm, there is no need of defining every edge case, a well defined agent should have a reasoning component capable of taking decisions even in unknown situations. The flexibility of an agent can be quantified across two main scales:

- reactive: delay required to respond to a change in the environment;
- proactive: ability of taking action in advance for the maximization of future goals.

Based on the environment and the agent itself there are different communication models but generally the agent perceives observations from the environment using sensors and performs actions against the environment employing actuators.

Another fundamental characteristic of an agent is its autonomy, the internal “brain” should handle the decision process with or without collected information, it should also evolve with respect to possible requirement changes.

Finally an agent can be designed to solve tasks or goals. When we talk about tasks we mean small objectives that must be completed to achieve the final bigger goal. On the other hand we have a goal agent, it takes the goal and it has to define autonomously a list of tasks to fulfill the assigned goal.

A. Multi-agent system

A multi-agent system, as suggested by the name, is a group of agents placed in same environment. There are two main type of interactions: cooperative and competitive. During this project we will focus on both of them in several stages of the development.

A competitive system is composed of many agents acting against each other, the goal can be divided into sub-goals, maximizing the personal reward and minimize the opponents’ ones.

A cooperative system is composed of many agents acting to maximize the shared reward, each agent must be capable of cooperate, coordinate and negotiate as much as possible.

In a cooperative system an important implementation detail is the communication, it should be fast and reliable with a the lowest possible delay.

B. Architecture

There exist multiple architectures that can be used to build an agent capable of acting in a given environment. We have decided to opt for the architecture described in the following pseudocode

Algorithm 1 Agent control loop

```

1: procedure AGENTCONTROLLOOP
2:    $B \leftarrow B_0$                                 ▷ Belief set initialization
3:    $I \leftarrow I_0$                                 ▷ Intention set initialization
4:   while true do
5:     perceive  $\rho$ 
6:      $B \leftarrow \text{update}(B, \rho)$                     ▷ Belief set update
7:      $D \leftarrow \text{options}(B)$                         ▷ Desires computation
8:      $I \leftarrow \text{filters}(B, D, I)$                   ▷ Intention update
9:      $\pi \leftarrow \text{plan}(B, I)$                       ▷ Plan computation
10:    while not (empty( $\pi$ ) or succeeded( $I, B$ ) or impossible( $I, B$ ))
11:      do
12:         $\alpha \leftarrow \text{hd}(\pi)$                       ▷ Get next action
13:        execute( $\alpha$ )                                ▷ Execute action
14:         $\pi \leftarrow \text{tail}(\pi)$                       ▷ Remove executed action
15:        perceive  $\rho$ 
16:         $B \leftarrow \text{update}(B, \rho)$                     ▷ Belief set update
17:        if reconsider( $I, B$ ) then
18:           $D \leftarrow \text{options}(B)$                     ▷ Desires computation
19:           $I \leftarrow \text{filters}(B, D, I)$                 ▷ Intention update
20:        end if
21:        if not sound( $\pi, I, B$ ) then
22:           $\pi \leftarrow \text{plan}(B, I)$                     ▷ Plan computation
23:        end if
24:      end while
25:    end while
26: end procedure

```

A belief set B is updated after every environment communication, it is then used to compute the desires set D which is filtered to create the intentions set I . Finally the intentions set is used in combination with the belief set to create a plan π to satisfy all the intentions. The plan is then executed action by action α . For every new information collected during the execution the belief set, the intentions set, and the desires set are updated. After the update a replan can be done, during this phase a new plan is computed.

II. DELIVEROO

Deliveroo is the environment provided to test the developed agents of the project, it is composed of a grid of tiles divided into two main categories: walkable (green and red) and not walkable (black). A walkable tile can be a delivery zone (red).

Several parcels, represented by cubes, are scattered around the map, the goal of the agent is to move across it collecting

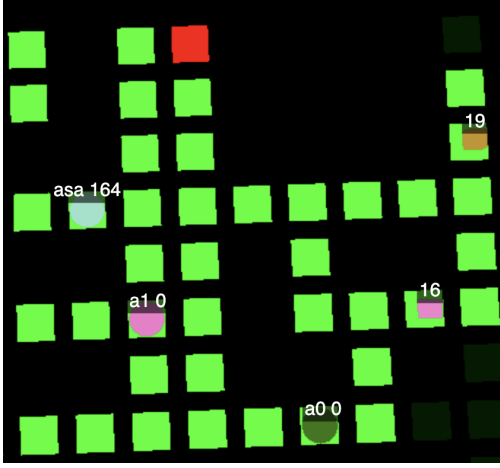


Figure 1. An example of grid map with different parcels and agents.

the largest number of parcels. The final objective is to deliver collected elements in one delivery zone.

The agent can move up, down, left, and right. It can pass through parcels and delivery zones but it is blocked by other agents that can play the game at the same time.

Multiple parcels can be carried by one agent, and the value of them continues to decrease

The agent can sense other agents and parcels, respectively within two well defined radiuses.

There are several parameters associated to the environment associated to various game elements:

- *parcel*: generation internal, reward decay, maximum number, and variance;
- *player movement*: steps number and speed;
- *agent*: amount and speed;
- *sensing*: parcel and agent sensing radius;
- *game clock*.

Given all the customizable parameters and the size of the map, the number of different combinations is huge. An agent must be well developed in order to be able to act in all of them, both alone or cooperating with other agents.

There are two main strategies that can be employed for the solution of the game:

- *single-agent*: the agent acts alone;
- *multi-agent*: the agent communicates with other mates to share information for the maximization of the summation of all the rewards;

III. PROBLEM ANALYSIS

Given the complexity of the environment several solutions have been implemented to overcome different difficulties and challenges.

A. Manhattan distance

The *Manhattan distance* is function that can provide a quick estimation of the distance between two points.

$$d_{ab} = |b_x - a_x| + |b_y - a_y|$$

In the deliveroo environment it is insufficient since many tiles can be non walkable, but at the same time it's fast heuristic to compute.

B. Problem parameters estimation

As explained in Section II, multiple parameters can be modified during the game initialization. Some of them are given to the agent, others are obscure and can be only estimated. We mainly focused on the estimation of player speed and parcel decay for a more accurate computation of the potential reward associated to a given parcel.

1) *Player speed*: Each time the agent moves the new position is communicated by the environment to the agent, we have decided to keep track of the timestamp in which the information is perceived. The list of timestamps can then be used to estimate the player velocity using the following algorithm:

Algorithm 2 Player speed estimation

```

1: procedure UPDATEMAINPLAYERSPEEDESTIMATION( $\mathcal{D}, s, \phi$ )
2:    $deltas \leftarrow []$   $\triangleright$  Initialize an empty array of deltas
3:   for each timestamp  $t_i \in \mathcal{D}$  do
4:      $deltas.append(t_i - t_{i-1})$   $\triangleright$  Delta of two consecutive timestamps
5:   end for
6:    $c \leftarrow s * (1 - \phi)$   $\triangleright$  Current speed contribution
7:    $n \leftarrow avg(deltas) * \phi$   $\triangleright$  New speed contribution
8:    $s \leftarrow c + n$   $\triangleright$  New estimation
9:   return  $c$ 
10: end procedure

```

The hyper-parameter ϕ is the learning rate that can be used to regulate the impact of the contribution with respect to the current speed estimation.

2) *Parcel decay*: Another important parameter for the correct definition of an agent is the parcel decay, it is the number velocity of the reward decrease. Similarly to the player speed estimation, the parcel decay estimation is computed from timestamp differences, where a timestamp is associated to an update of a visible parcel.

Algorithm 3 Get parcel decay estimation

```

1: procedure GETPARCELDECAYESTIMATION( $\mathcal{D}$ )
2:    $deltas \leftarrow []$   $\triangleright$  Initialize an empty array of deltas
3:   for each timestamp  $t_i \in \mathcal{D}$  do
4:      $deltas.append(t_i - t_{i-1})$   $\triangleright$  Delta of two consecutive timestamps
5:   end for
6:   return  $deltas$ 
7: end procedure

```

Algorithm 4 Parcels decay estimation

```

1: procedure UPDATEPARCELSDECAYESTIMATION( $\mathcal{P}, d, \phi_2$ )
2:    $deltas \leftarrow []$   $\triangleright$  Initialize an empty array of deltas
3:   for each parcel  $p_i \in \mathcal{P}$  do
4:      $deltas.concat(getParcelDecayEstimation(p.timestamp))$ 
5:   end for
6:    $c \leftarrow d * (1 - \phi_2)$   $\triangleright$  Current decay contribution
7:    $n \leftarrow avg(deltas) * \phi_2$   $\triangleright$  New decay contribution
8:    $d \leftarrow c + n$   $\triangleright$  New estimation
9:   return  $d$ 
10: end procedure

```

The learning rate ϕ_2 can be used to regulate the impact of the contribution with respect to the current parcel decay estimation.

C. Probabilistic model

In the environment there may be multiples competitive agents, their ability of picking up parcels highly influences the value of a parcel. For this reason we have devised a penalty value based on a probabilistic model capable of taking into consideration the possible opponents' plans.

The main idea behind the probabilistic model is the following assertion: if there is a parcel and I am the closest agent I can reach it faster than any other agents, consequentially that parcel should be taken more into consideration, even if its value is lower than other further parcels.

This assertion can be modeled as follow

$$\text{penalty probability} = \frac{\sum_{a \in \mathcal{A}} \frac{d_{max} - d_{pa}}{d_{max}}}{|\mathcal{A}|}$$

with \mathcal{A} the set of opponent agents, d_{max} the maximum distance between the parcel and the union between opponents agents, main player, and cooperative agents, d_{pa} the distance parcel opponent agent.

D. Potential parcel score

The decision process behind the choice of a parcel is one of the key elements for the definition of a good agent. Many elements and metrics have been taken into consideration to better estimate the potential gain of a parcel. The final reward is computed as:

$$r_f = r - \left(d_{ap} * \frac{s_a}{decay} \right) - \left(d_{min} * \frac{s_a}{decay} \right) - r * \text{penalty probability}$$

where d_{ap} is the distance between the agent and the parcel, s_a is the estimated speed, d_{min} is the minimum distance between the parcel and any delivery zone, and penalty probability is the probability computed on Section III-C.

The resulting formula takes into reward the lost of approaching the parcel and delivery it to the closest delivery zone, moreover the probabilistic model can provide a rough estimation of other agents' intentions.

E. Distances cache

A distance cache was maintained to save some computational power, every time a plan is computed the distance between the starting point and any other tile in the path is saved in the cache.

The cache is then used in many parts of the code.

F. Replan

Qui non so se per ora il replan dopo 5 tentativi sia interessante da menzionare.

IV. SOLUTIONS

Qui spieghiamo le varie soluzioni presentate.

Forse questo capitolo e' ridondante, non saprei.

V. BENCHMARKING

VI. CONCLUSION

VII. REFERENCES