

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 to 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 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 procedure

```

A belief set B is updated after every environment communication, and 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 execute action by action α . Once the agent senses new information from the environment, the belief set, the intentions set, and the desires set are updated. After this update, the agent can decide whether to generate a new plan or stick with the current one.

II. DELIVEROO

Deliveroo is the playground environment which has been provided to develop our BDI agents. It consists of a delivery-based game, where agents move on a bidimensional grid. The objective of the game is to pick up parcels that are scattered across the map and deliver them as fast as possible in the tile classified as delivery zones. Each parcel has an assigned reward value, which can decay over time, that is assigned to the agent that delivers it.

The game can be played by either a single agent or multiple agents too. In the latter case, agents are solid entities and can block the road to each other.

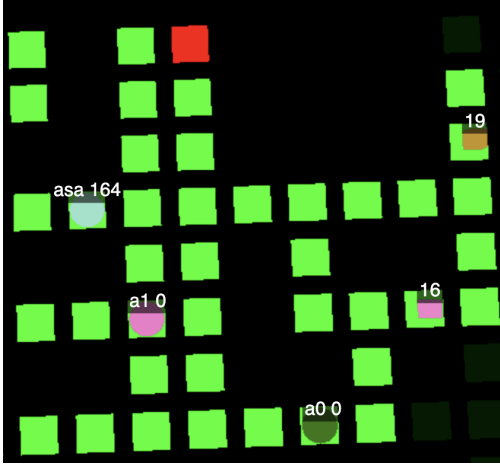


Figure 1. An example of a Deliveroo map.

Agents usually don't have the full perception of the world, but can perceive only other agents and parcels within a limited distance radius. Agents can also carry and deliver multiples parcels at the same time.

Deliveroo has been designed to offer many different game scenarios, thanks to the usage of many parameters: - *parcel*: generation interval, reward decay, reward value distribution; - *player*: steps number and movement speed; - *external agent*: quantity and movement speed; - *sensing*: parcel and agent sensing radius; - *map*: definition of the grid board and tile classification (walkable, non-walkable, delivery) - *game clock*.

Specifically, this exam project is split in two deliverables, each one with a different solution strategy: - *single-agent*: the agent acts alone trying to maximize its score; - *multi-agent*: the agent communicates with others to share information and to maximize the overall score of the group.

III. SINGLE AGENT IMPLEMENTATION

Given the complexity of the environment, several ad-hoc 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.

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

Even though this function might be insufficient in the Deliveroo game since maps can have non-walkable tiles, it is still a decent and fast to compute heuristic that can be used to approximate distances.

B. Problem parameters estimation

As explained in Section II, multiple parameters can be modified during the game initialization. While some of them are explicitly communicated to the agent, others are obscure and can be only estimated. For this reason, in order to build a more precise player reward estimation, we developed a learning mechanism to approximate the player speed and parcel reward decay over time.

1) *Player speed*: When our agent moves in the map, it is able to track its position over time: in this way, by using the historical positions alongside with timestamps the agent can estimate its own speed by 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

```

Where ϕ is the learning rate hyper-parameter that can be used to regulate the impact of the last measured instant speed with respect to its historical value.

2) *Parcel decay*: In addition to the agent speed estimation, our agents are able to estimate the parcel reward decay over time. Similarly to the player speed estimation, the parcel decay is computed from timestamp differences, where each timestamp is associated to a sensed reward 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:      $deltas.concat(getParcelDecayEstimation(p_i.timestamps))$ 
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

```

Where ϕ_2 is the learning rate and can be used to regulate the contribution of the past estimations with respect to the current estimated parcel decay.

C. Probabilistic model

In the environment there may be multiple competitive agents, and 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". More formally:

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

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 parcel selection is one of the key elements when defining a good agent. Many elements and metrics have been taken into consideration to better estimate the potential reward gain of a parcel. More formally, our agents compute the final reward 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 the closest delivery zone, and penalty probability is the probability computed on Section III-C.

The resulting formula takes into account multiple factors: - the agent reward that is spent to approach the parcel and delivery it to the closest delivery zone (minimum cost to deliver that package) - rough estimation of other agents' intentions through their distance from parcels by using a probabilistic model

E. Distances cache

To save computation power and to provide a more precise reward estimation, a cache is maintained to store distances between tiles across the map. Every time a plan is generated, the distance between the starting point and any other tile in the path is stored in the cache. Since this approach does not guarantee that the generated path is the shortest route between two tiles, cache entries are updated once a smaller value is found. Therefore the caching approach is meant to improve over time as a sort of learning mechanism. In general, the cache is used many times in the codebase: in case of a cache miss, the agent uses the Manhattan distance as fallback.

F. Replan

Since Deliveroo is a dynamic game with potentially multiple agents acting at the same time, we implemented a mechanism to replan the current agent's actions when it fails to perform a move for a specific amount of time. This functionality enables agents to avoid them from sticking in narrow and crowded areas of the map.

IV. MULTI AGENT IMPLEMENTATION

The communication protocol is based on the provided library that offers multiple endpoints to handle different messages: say, shout, ask, and broadcast.

A. Information sharing agents

During the game, teams can share the data they sense from the environment between each other. For simplicity, for this specific purpose our agents use broadcast messages that are visible to every agent that is connected to the game.

In this way, every agent can build its one plan based on the information which is sensed by all the other agent across the map.

Moreover, agents also communicate by broadcast message the packages that they are going to pickup with the current plan they just generated. This mechanism helps in reducing the number of collisions of agents in the map and avoiding useless longer paths to pickup already taken parcels.

B. Leader-members agents

The leader negotiation is a fundamental step in the multi-agent architecture where one agent is chosen to be the leader who computes all the plans for the other agents.

1) *Leader negotiation*: The election of the leader is a well known problem in the computer science literature, nevertheless we have decided to keep a simple solution since the focus on the project was something else. The negotiation is achieved thanks to two message types: askforleader and leader. The former is sent by every agent as a broadcast communication when they connect to the game (when they receive for the first time their position): as the name says, the message simply asks if a leader has already been elected. The latter is used to communicate the leader identity and it is used only by the agent who has been elected as the leader.

After the askforleader message is sent, a timeout of 2.5s is set. If no answer by an actual leader is received during that time interval, it means that the game has no already elected leader and consequentially the agent who asked for the leader identity will be elected as the leader itself. At this point, this agent communicates as a broadcast message that it has just elected itself as the leader.

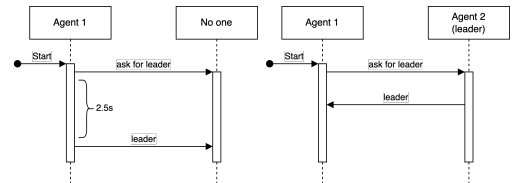


Figure 2. Leader negotiation.

Please, consider that with this simple implementation very rare scenarios where two (or more) agents exactly connect at the same time could generate a race condition that would likely create more than one active leader.

V. BENCHMARKING

VI. CONCLUSION

VII. REFERENCES