Using POMDP with raw observations for detecting and recognizing objects of interest

Jean-Loup Farges¹, Guillaume Infantes¹, Charles Lesire¹ and Augustin Manecy¹

Abstract—Using POMDP-based decision-making is an appealing choice for robotic missions, especially because this method explicitly model uncertainties related to sensors. However, it is still hard to obtain successful experiments when sensors provide noisy or ambiguous data. We think one of the main reasons is that observation models are most of the time poorly reliable, especially when using vision-based sensors and processings. In this paper, we propose an architecture that, instead of using an external classifier, connects directly the outputs of detectors to the decision making component. We propose an adaptation of POMDPs to directly use multi-dimensional continuous data delivered by detectors or image processing algorithms. This adapted POMDP framework outperforms the classical approach tested experimentally with an autonomous aerial vehicle detecting and recognizing objects of interest.

I. INTRODUCTION

Real-world autonomous robotic missions have to deal with uncertainty because sensor outputs are noisy and not always reliable. A popular framework for dealing with uncertainties in sequential decision making is Partially Observable Markov Decision Processes (POMDP) [1] because they model sensor noise explicitly with observation probabilities, along with system state estimation. This framework is useful for several kinds of missions including manipulation [2], navigation [3], or high level tasks [4]. Classical POMDP approaches used in robotic detection and recognition applications often define the observation model by discretizing the raw detector outputs as observation symbols using classification. Such traditional architecture for integrating image processing and decision-making based on POMDP is illustrated as Fig. 1. It is for instance used in [5] with a simple classifier, based on the comparison of the highest score with a threshold, to detect and recognize vehicles while exploring an area. Such approaches could be improved using an advanced classifier, learned by kernel-based methods [6] or deep learning [7]. While the latter is the new state of the art in discrete classification, it needs very large hand-tailored augmented datasets in order for the optimization algorithm to achieve the very good performances reported. This can be a limiting factor in very specific setups where extensive data collection is hard, typically in robotic experiments.

The problem of using discrete observations is that classifiers introduce over simplistic observation models within the decision component: using discrete observations and a confusion matrix eases the definition of the Bayes rule for belief update, but loses information about the nature of the several continuous observations. Hoey and Poupart have



Fig. 1: Traditional approach: detectors scores are used by a classifier to identify the most probable class of the object.

indeed shown, in simulation on a two-states two-decisions POMDP model, that considering continuous observation variables instead of the result of a classification process improves the quality of the policy [8].

In this paper, we propose to use the architecture illustrated as Fig. 2. In this architecture, several detector scores are directly used as an input of the decision component, without any classifier. In such setup, we can incorporate a range of detectors not necessarily linked to the number of target classes and directly take decisions based on these observations.



Fig. 2: Proposed approach: detector scores are directly used by the decision module.

Trying to make a step towards the use of deterministic policies for discrete state and multiple continuous observations in field robotics, we have implemented our approach on

 $^{^1} authors are with ONERA – The French Aerospace Lab, Toulouse, France firstname.lastname@onera.fr$

an object recognition task by an autonomous aerial vehicle. Then we have compared to the standard POMDP approach. This paper brings two contributions: a POMDP method for integrating multiple continuous observations of discrete states, and its implementation and evaluation on a real experiment. Next section introduces the necessary background on POMDP, and the related works on considering multiple continuous observations. Section III then presents the POMDP model with continuous multi-dimensional observations suited to an object recognition task. An algorithm to solve such a POMDP is proposed in Section IV. Section V then presents the experimental setup, and results are discussed in Section VI.

II. BACKGROUND AND RELATED WORKS

A POMDP is classically defined by [1]:

- S a finite set of states of the system;
- \mathcal{A} a finite set of actions;
- T: S×A → Π(S) the state-transition function, giving for each state and robot action, a probability distribution over next states; we write T(s, a, s') for the probability of ending in state s', given that the robot takes action a in state s;
- R: S×A → ℝ the reward function, giving the expected immediate reward gained by the robot for taking each action in each state; we write R(s, a) for the expected reward for taking action a in state s;
- Ω a finite set of observations the robot can experience after doing an action;
- O: S × A → Π(Ω) the observation function, which gives, for each action and resulting state, a probability distribution over possible observations; we write O(s', a, o) for the probability of making observation o given that the robot took action a, resulting in state s'.

Controlling a POMDP is done by computing a policy π that maximizes the expected total reward. As the state of the world is not directly observable, a general scheme for computing π is to distinguish two processes: (1) the estimation (or filtering) of the state of the world as a *belief state*, and (2) the computation of the policy over the belief states. Historically, algorithms for solving POMDPs (i.e., computing π) consider that S, A and Ω are discrete and finite, as in previous definition.

In order to be able to use all the information available from the detectors within the decision component (see Fig. 2), the POMDP model should take into account multi-dimensional and continuous observations. The set of observations can then be defined as $\Omega = \mathbb{R}^Z$ where Z is the number of detectors.

First works toward considering continuous spaces extended the POMDP framework to continuous state space [9] and continuous action space [10]. Algorithms to solve these POMDPs are known as point-based algorithms; they use sampling to explore the continuous spaces.

Other works have extended POMDPs to continuous observations. [11] extends a point-based algorithm to continuous state space and continuous observation space. The authors do not explicitly update the belief in a continuous multidimensional manner, but instead, sample the belief space in order to consider a finite planning graph for policy computation. A scheme for discretizing observations is presented in [8]; again it does not take into account full continuous multi-dimensional observation probabilities as we propose. [12] considers continuous state, action and observation spaces, but restrained to a Gaussian belief representation, which leads to an update using an extended Kalman filter. Such approach, inspired from Optimal Control theory, is less generic than the one we propose, as we do not use such restrictive hypothesis for belief representation.

Some approaches directly learn parametric policies as functions of the histories of observations without explicitly computing the belief state (e.g., [13] uses Deep Neural Networks). Such setup can naturally be used with continuous state and observation spaces but cannot take into account an a priori belief. Moreover, we propose an approach that considers as much information as possible from the observations models in a founded way. A POMDP model with a continuous observation space has been proposed in [14], and applied to a scenario where several robots have a noisy observation of their own position. Continuous observations are used (using a regression algorithm) to probabilistically select the decision state, which relies in a discrete space, like in other finite-state automaton-based methods. These observations are not used to update a continuous belief state, as we propose to do.

To the best of our knowledge there is no work computing a deterministic policy from a POMDP defined by a discrete state space and a continuous multi-dimensional observation space. Moreover, none of the aforementioned works in robotics have been implemented on an experiment with a real robot.

III. POMDP MODEL FOR OBJECT RECOGNITION

A. Epistemic MDP

In this paper, we consider a specific POMDP model known as the Epistemic Markov Decision Process (EMDP) [15]. In an EMDP, the state of the world does not change: the actions performed by the robot are aimed at improving its knowledge about the world state, through new observations. These observations are used to update the belief state. An EMDP is then a special case of a POMDP as defined in Section II, with the state-transition function T defined as:

$$\forall s, s' \in \mathcal{S}, \forall a \in \mathcal{A}, \ T(s', a, s) = \begin{cases} 1 & \text{if } s = s' \\ 0 & \text{if } s \neq s' \end{cases}$$
(1)

An EMDP classically has a finite horizon $H \in \mathbb{N}$. The reward function R is then defined in order to represent the epistemic quality of the final belief.

The EMDP model is well suited to object detection and recognition missions, as it models that the robot tends to improve its knowledge about the true state of the world without modifying this state. In the following, we then focus our contribution on solving EMDP with multi-dimensional continuous observations. Note that this focusing do not restrict the contribution of this paper. First, [15] has shown that solving an EMDP is as complex as solving a classical POMDP. Second, in the case where the mission is not purely epistemic, our approach may be augemented with a state filtering step in a straightforward way. We would like to emphasize that for some missions, it may be interesting to completely separate the epistemic objectives from other mission goals, in order to interleave several purely knowledgegathering steps with more classical acting steps.

B. Object recognition model

The object recognition problem can then be modelled as an EMDP with stochastic continuous multi-dimensional observations. We consider a set of M classes $C = \{c_1, \ldots, c_M\}$ of objects of interest that are distributed in a set of Nlocations $\mathcal{L} = \{l_1, \ldots, l_N\}$. To simplify the description of the approach, we consider that each location contains at most one object, and that a specific class represents the fact that a location is empty. Consequently, the set of states is defined as $S = C^N$. For a state $s = (s_1, \ldots, s_N)$, $s_i \in C$ represents the class of the object in location l_i .

The set of actions \mathcal{A} contains one action for observing each location: $\mathcal{A} = \{a_1, \ldots, a_N\}$. After executing each action, the robot will observe the state through its sensor, and a score is computed by each detector. As illustrated in Fig. 2, each observation is then a vector of continuous scores. When action a_i is performed, it results in an observation $o_i \in \Omega$ giving observation about the object located in l_i . o_i is a vector $o_i = (o_{i,1}, \ldots, o_{i,Z})$ where each $o_{i,j}$ is the score obtained from the *j*-th detector.

The observation function O of an EMDP is defined as

$$O(s, a, o) = p(o|s, a) \tag{2}$$

In our model, as each action a_i corresponds to observing a sub-part s_i of the state (corresponding to location l_i) through observation o_i , we can note the observation function as:

$$\forall i, \ O(s_i, a_i, o_i) = p(o_i | s_i, a_i) \tag{3}$$

Knowing that each s_i contains an object of a class $c_j \in C$, O is then defined by a family of probability density functions $p(o_i|a_i, s_i = c_j)$, that can be shorten into $p(o_i|c_j)$.

IV. DYNAMIC PROGRAMMING ALGORITHM

Solving (PO)MDPs consists in computing an action policy that optimizes some numeric criterion V, named value function. In this section, we first describe the belief space we consider (IV-A), and then we propose a non-classical criterion that suits the model previously defined (IV-B). The core of the POMDP resolution mechanism is then adapted by formulating, in IV-C, a Dynamic Programming equation for our model and using a sampling scheme over the possible observations for belief update in IV-D. We finally give a synthetic view of the whole algorithm in IV-E.

A. Belief space

In the object recognition model, each observation gives information about only one location. Moreover, the classes of the several objects are independent. We thus decompose the belief state B into a tuple of beliefs b_i , $i \in [\![1, N]\!]$, and each b_i is a vector $b_i = (b_i(c_1), \ldots, b_i(c_M))^T$ giving the estimated probabilities that location l_i is occupied by an object of class c_i , $j \in [\![1, M]\!]$.

Then, for any $c_j \in C$, the Bayes rule gives:

$$\forall i, k \in \llbracket 1, N \rrbracket, i \neq k, \quad b_i(c_j | o_k) = b_i(c_j) \tag{4}$$
$$\forall i \in \llbracket 1, N \rrbracket, \quad b_i(c_j | o_i) = \frac{b_i(c_j)p(o_i | c_j)}{\sum_{c' \in \mathcal{C}} b_i(c')p(o_i | c')} \tag{5}$$

The vector form of Eq. (5) can be written as:

$$b_i(|o_i) = b_i \odot m(b_i, o_i) \tag{6}$$

with \odot the component wise multiplication of vectors and

$$m(b_i, o_i) = \frac{1}{\sum_{c' \in \mathcal{C}} b_i(c') p(o_i | c')} \cdot (p(o_i | c_1), \dots, p(o_i | c_M))^T \quad (7)$$

B. Criterion representation

As defined in III-A, EMDPs classicaly define the reward function so that it represents the quality of the final belief. The optimal value of the criterion, V^* , is then defined as being the reward associated to this final belief. Taking into account the factorization of states and beliefs of different locations, we define the optimal value of the criterion V^* with respect to the current belief as the minimum of the expected average distance between the belief and the true state of the world. Let d be the distance between two classes:

$$\forall c, c' \in \mathcal{C}, \ d(c, c') = \begin{cases} 1 & \text{if } c \neq c' \\ 0 & \text{if } c = c' \end{cases}$$
(8)

This criterion can be written as:

$$V^*(B) = \min_{s \in \mathcal{S}} \frac{1}{N} \sum_{i=1}^N b_i \cdot d(s_i)$$
(9)

with $B = (b_1, ..., b_N)$ and $d(s_i) = (d(s_i, c_j))_{1 \le j \le M}$.

C. Dynamic Programming principle

In the following, we use a dynamic programming scheme, i.e. we will define a sequential problem over H timesteps. $a^{(k)}$, $o^{(k)}$ are the action chosen and the observation done at timestep k; the belief at timestep k is $B^{(k)} = (b_1^{(k)}, \ldots, b_N^{(k)})$. The value function is classically represented by as a set of $L^{(k)}$ alpha-vectors $\alpha_j^{(k)}$, $j \in [\![1, L^{(k)}]\!]$. Then:

$$V^*(B^{(k)}) = \min_{j \in [\![1, L^{(k)}]\!]} B^{(k)} \cdot \alpha_j^{(k)}$$
(10)

The representation of the criterion above can be matched with alpha-vectors representation, and so we have, at decision horizon H:

$$V^*(B^{(H)}) = \min_{j \in [\![1, L^{(k)}]\!]} B^{(H)} \cdot \alpha_j^{(H)}$$
(11)

We now have to define how to compute recursively the $V^*(B^{(k)})$, i.e. define the Bellman backups for our problem.

D. Bellman backup

Let $E_{o^{(k)}}$ denote the mathematical expectation with respect to $o^{(k)}$. Bellman's optimality principle then allows to state a relation between $V^*(B^{(k)})$ and $V^*(B^{(k+1)})$:

$$V^*(B^{(k)}) = \min_{a^{(k)} \in A} E_{o^{(k)}}(V^*(B^{(k+1)}(|a^{(k)}, o^{(k)}))) \quad (12)$$

where

$$b_i^{(k+1)}(|a^{(k)}, o^{(k)}) = \begin{cases} b_i^{(k)} & \text{if } a^{(k)} \neq a_i \\ b_i^{(k)} \odot m(b_i^{(k)}, o^{(k)}) & \text{otherwise} \end{cases}$$
(13)

Thus assuming that $V^*(B^{(k+1)})$ is defined by a set of alpha vectors $\alpha_j^{(k+1)}$ with $j \in [\![1, L^{(k+1)}]\!]$:

$$V^{*}(B^{(k)}) = \min_{a^{(k)} \in A} E_{o^{(k)}} \min_{j \in [\![1, L^{(k+1)}]\!]} \sum_{i=1}^{N} (b_{i}^{(k)} \odot m(b_{i}^{(k)}, o^{(k)}) \times \alpha_{j,i}^{(k+1)})$$
(14)

Approximating the mathematical expectation with P random samples leads to:

$$V^{*}(B^{(k)}) \cong \min_{a^{(k)} \in A} \frac{1}{P} \sum_{i=1}^{P} (\min_{j \in [\![1, L^{(k+1)}]\!]} (B^{(k)} \odot m(B^{(k)}, o_{i})) \times \alpha_{j}^{(k+1)})$$

$$m(B^{(k)}, o_{i})) \times \alpha_{j}^{(k+1)})$$
(15)

where o_i is randomly sampled taking into account belief $B^{(k)}$ and action $a^{(k)}$. For given belief, action and sample, the argument of the second minimum is $j^*(B^{(k)}, a^{(k)}, i)$. Moreover for a given belief the argument of the first minimum is $a_*^{(k)} = a^*(B^{(k)})$. Thus:

$$V^{*}(B^{(k)}) \cong B^{(k)} \times \left(\frac{1}{P} \sum_{i} \alpha_{j^{*}(B^{(k)}, a_{*}^{(k)}, i)}^{(k+1)} \odot \right)$$

$$m(B^{(k)}, a_{*}^{(k)}, o_{i}))$$
(16)

The value $\alpha^{(k)}(B^{(k)})$ for the neighborhood of $B^{(k)}$ is then:

$$\alpha^{(k)}(B^{(k)}) \cong \frac{1}{P} \sum_{i} \alpha^{(k+1)}_{j^*(B^{(k)}, a^{(k)}_*, i)} \odot m(B^{(k)}, a^{(k)}_*, o^i)$$
(17)

If the result of the product of this $\alpha^{(k)}(B^{(k)})$ by $B^{(k)}$ is lower than the result of the product of other $\alpha^{(k)}(B'^{(k)})$ by $B^{(k)}$, this $\alpha^{(k)}(B^{(k)})$ may contribute to the definition of the function $V^*(B^{(k)})$ and is to be considered for inclusion in the set of alpha vectors $\alpha_j^{(k)}$ that define $V^*(B^{(k)})$.

E. Policy computation algorithm

In order to compute the policy, we have implemented a simple uniform sampling over the whole belief state space, and applied the DP computations over the grid points, for simplicity of implementation reasons. For better performance (in terms of scalability), our contribution should be used with smarter picewise-linear functions implementation like basis functions [16], and principled sampling as in SARSOP [17].

This gives, for a belief state $B^{(k)}$, for each action $a_i^{(k)}$:

- to sample possible observations;
- to update the belief wrt. current action $a_i^{(k)}$ and sampled observation using Eq. (5) and (13)
- to compute the corresponding alpha vector and expected cost using Eq. (17) and (15);
- to compute average (over observations samples) of alpha vectors and expected costs.

When this is done for each action, we choose the best action (in term of expected cost) and add the corresponding (averaged) alpha vector to the set of alpha-vectors.

V. EXPERIMENTAL SETUP

In this section, we first briefly describe the mission performed by the drone, then we present the vision-based algorithm we use as detectors, and how we have defined the corresponding observation model. We also introduce a classical POMDP approach for solving this problem, to which we will compare.

A. Mission and platform

Experiments have been made using a Parrot AR.Drone with a downward camera. The drone is localized using a motion capture system, and is controlled by a PID controller, to which the decision system can send position and heading targets. The possible classes of the searched objects are: *empty*, *jacket* and *box*. These objects can be located in any of the two possible locations l_1 and l_2 . A typical setup is shown in Fig. 3.



Fig. 3: In this scenario, the AR.Drone has to find a yellow jacket in location l_1 and an orange box in location l_2 .

The policy is computed using the approach proposed in this paper, given an initial belief about occupancies of the several locations. The policy then indicates, for each belief, which action to perform. When performing an action a_i to observe location l_i , the control architecture of the drone randomly chooses a pose to observe the corresponding location. These possible poses correspond to observing the location from the north, south, east, and west orientation, at two possible altitudes.¹

¹This choice has been made to have a great variability in the observations.



Fig. 4: Description of the pipeline of the two detectors detecting a box and a jacket.

Once the action is finished, the observation arrives as an array of scores coming from the several detectors. The belief is then updated using Eq. (5) for the current location, and the policy gives the next action to perform according to the current belief. At the end of the last step, the most likely occupancy state for each location is returned.

B. Detectors

We consider here that the robot is endowed with two detectors, dedicated to find a specific object. The first detector has to find box objects while the second has to find *jacket* objects. Each detector consists of an object tracker implemented as a vision-based algorithm which processes images provided by the downward camera (the sensor). Each object tracker delivers a score indicating the probability that an object of the corresponding class is present in the current image. We have implemented an algorithm based on the well known CAMSHIFT algorithm [18]. Details about our implementation can be found in [19]. This algorithm looks for an object into an image according to its color signature (e.g., a reference histogram of the object). Figure 4 illustrates the complete pipeline of this algorithm: RGB images are first converted to HSV images to easily extract the hue component. Thresholds on the saturation and value component are applied to filter not reliable pixels. Then the hue image and a reference histogram are used to compute a backprojection using unmasked pixels. The CAMSHIFT algorithm then provides a bounding box containing the object for which we determine the corresponding histogram. In order to deliver the scores o_i corresponding to each detector, we finally compute the cross-correlation between the reference histogram of the object and the histogram of the bounding box.

By correctly defining the reference histograms and the thresholds, this algorithm would give very good results for

our experiment. The objects are indeed quite different, and there is no disturbing object in the environment. In order to have a more interesting situation, we decided to not spend a lot of time in defining the histograms and thresholds. It results in more noisy detectors, which are more representative of a real scenario in which the objects to recognize are more ambiguous. Figure 4 shows a situation where an orange box is detected in the shadow of the jacket, with a score of 0.395. Other noisy situations are shown in Fig. 5c to 5d.



(c) box: 0.466; jacket: 0.769

(d) box: 0.338; jacket: 0.025

Fig. 5: Some observations made during experiments

C. Observation model

As described in Sect. III, the observation model is defined by a family of probability density functions $p(o_i|c_j)$, for $i \in [\![1, N]\!]$ and $j \in [\![1, M]\!]$. In the experiment, we have decided to use a representation with parameters depending on the class c_j but not on the location l_i . From the detection algorithm described in the previous section, it seems natural to consider that the score of the detectors will only depend on the class of the object, i.e. on its reference histogram, independently from the locations.

We define the observation model as a multivariate normal distributions $p(o_i|c_j) = f_Z(\mathcal{N}_Z(\mu_{c_j}, \Phi_{c_j}))$, where μ_{c_j} and Φ_{c_j} are respectively of dimension Z (the number of detectors) and $Z \times Z$ and we saturated it (in order to respect the bounded nature of scores) by a sigmoid function $f_Z(x_1, ..., x_Z) = (f(x_1), ..., f(x_Z))$ and $f(x) = \frac{1}{1+e^{-x}}$. This model is then identified from the scores of a data set.

In our experiment, there are three possible classes *empty*, *jacket* and *box*, and two detectors, producing two scores: the *jacket score* and the *box score*. Figure 6 plots the scores obtained on the data set in order to identify the model prior to computing the policy. This data set has been obtained by placing a jacket, a box, or no object in the experiment environment, and performing observations using several conditions (position, orientation, time of day, luminosity). The data set contains a hundred of pairs of scores for each class.



Fig. 6: Scores measured for three object classes

For each class c, μ_c and Φ_c are computed as the mean and covariance matrix of these scores. These resulting distributions are illustrated in Fig. 7. One can see that these model are not well separated, especially as the output from the sensors are between 0 and 1, which is where the gaussian overlap the most. In our opinion, this justifies to use the well-separated beliefs in order to take into account the bidimensional structure of the observation.

D. Initial belief

In order to evaluate the impact of the initial belief on the quality of the computed policy, we have defined several initial beliefs:

- 1) an equiprobable belief: for each location, the same probability of $\frac{1}{3}$ is given to each class;
- an almost empty belief: for each location, initial belief of the empty class is 0.8, other classes are 0.1 each;
- 3) an asymmetrical belief: for l_1 , jacket and box have the same probability of 0.4, empty has 0.2; for l_2 , box has 0.8, others have 0.1 each.

The rationale of these initial beliefs is to represent: (1) no a priori knowledge, (2) most of the time locations are empty,



Fig. 7: Identified normal distributions for the 3 classes

(3) some information coming from a human-begin (e.g., an operator analysing a high altitude image of the locations).

E. Classical POMDP setup

We have compared our approach to a more classical POMDP, where the scores from the detectors are first used by a classifier to identify the most probable class of the observed object (see Fig. 1). We have then first implemented a classifier using a Gaussian kernel trained with LibSVM [20]. The observation model learned from the data set described in the previous section is given in Table I.

TABLE I: Discrete observation model probabilities

	P(O = jacket / T)	P(O = box / T)	P(O = empty / T)
T = jacket	0.990	0.010	0
T = box	0.293	0.707	0
T = empty	0.505	0.316	0.179

In this model, there is a very strong bias towards having the *jacket* observation symbol. On the contrary, the *empty* observation does not happen often, but is a very good clue about the ground truth.

We have then used the PBVI algorithm [21] to compute the policy for a classical discrete POMDP using the observation function defined by Tab. I. Indeed, we think that suboptimality induced by sampling in PBVI or SARSOP [17] are negligible in our experimental setup, more precisely suboptimality is only caused by the loss of information induced by the observation model.

VI. RESULTS

We have defined several scenarios in order to evaluate our approach. These scenarios are defined by:

- the real *state* (i.e. which object is in which location);
- the initial belief as defined in Sect. V-D;
- the architecture: classical POMDP (classifier plus PBVI), or our approach integrating multi-dimensional continuous observations within the POMDP model (noted raw_obs in the following figure).

The experimented *states* are: (1) both locations are empty, (2) a box in l_1 , (3) a box in l_2 , (4) one jacket in each location, and (5) a jacket in zone l_1 and a box in l_2 . We then performed 10 runs for each scenario, leading to a total of 300 experiments (5 states, 3 initial beliefs, 2 policies, 10 runs). We compare the performance of our approach with respect to the classical POMDP architecture by comparing the rate of correct classifications. A classification is correct if the class of an object corresponds to the most probable class according to the final belief.

The results obtained from the 300 experiment are a correct classification rate of 0.5117 for our proposed method, versus 0.4767 for the classical approach. The improvement raised by our approach is then of 7%.

These classification rates are however quite low, indicating that both methods faced perception situations that did not correspond quite well to the models. We have then analysed the scores observed during the experiments, presented in Fig. 8. The repartition of the scores is very different from the scores used to identify the POMDP models (Fig. 6), i.e. identified normal distributions (as shown in Fig. 7) would have significantly difference parameters. These differences may be due to some illumination conditions not captured during the creation of the data set, or by a degradation of the quality of the physical sensor (the drone camera).



Fig. 8: Scores observed during the experiment

Figure 9 shows the scatter plot of the correct classification rates for the two policies, according to the several states. Plots above the first diagonal have better scores with our policy using raw observations, while plots below the diagonal have better scores with the classical approach.

The classical POMDP approach mainly outperforms our approach in the states containing two yellow jackets. This results come from the learning of the classifier (Section V-E), that has a strong bias in classifying an object as a jacket than any other class. For other states, our approach globally outperforms the classical approach.

VII. CONCLUSION

In this paper, we have proposed a new POMDP-based architecture for detecting and recognizing objects of interest,



Fig. 9: Correct classification rates of the two methods for several scenarios.

which is a common task in robotics. Our architecture directly integrates multi-dimensional continuous scores coming from several detectors into the POMDP model. It does not rely on a classifier that aggregates the several detector scores into one symbolic observation, but uses the raw observation directly when updating the belief state and computing an optimal policy. We have defined our POMDP model as an Epistemic MDP with continuous and multi-dimensional observations, and we have proposed a Dynamic Programming algorithm to solve this model.

Finally, we have evaluated our approach by comparing to a classical POMDP architecture on an object detection mission by an UAV, in a large set of situations. These experiments have shown that our approach outperforms the classical POMDP approach.

Future work will consist in evaluating our method on a wider range of situations, including different object classes, and more complex missions. We will also investigate other POMDP models that may provide less precise but more robust results in order to face situations were the observed scores may be more ambiguous. For instance, extending possibilistic POMDP [22] with continuous multi-dimensional observations seems a promising approach.

REFERENCES

- L. Kaelbling, M. Littman, and A. Cassandra, "Planning and acting in partially observable stochastic domains," *Artificial Intelligence*, 1998.
- [2] J. Pajarinen and V. Kyrki, "Robotic manipulation of multiple objects as a POMDP," Artificial Intelligence, 2015.
- [3] J. Pineau and G. J. Gordon, "POMDP Planning for Robust Robot Control," in ISRR, 2007.
- [4] J. Pineau and S. Thrun, "High-level robot behavior control using POMDPs," in AAAI, 2002.
- [5] C. P. C. Chanel, F. Teichteil-Königsbuch, and C. Lesire, "Multi-Target Detection and Recognition by UAVs Using Online POMDPs," in AAAI, 2013.
- [6] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel Methods in Machine Learning," *The Annals of Statistics*, vol. 36, no. 3, 2008.
- [7] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, http://www.deeplearningbook.org.

- [8] J. Hoey and P. Poupart, "Solving POMDPs with continuous or large discrete observation spaces," in *IJCAI*, 2005.
- [9] J. M. Porta, N. Vlassis, M. T. Spaan, and P. Poupart, "Point-Based Value Iteration for Continuous POMDPs," *JMLR*, 2006.
- [10] M. T. Spaan and N. Vlassis, "Planning with Continuous Actions in Partially Observable Environments," in *ICRA*, 2005.
- [11] H. Bai, D. Hsu, and W. Lee, "Integrated perception and planning in the continuous space: A POMDP approach," *IJRR*, vol. 33, no. 9, 2014.
- [12] J. Van Den Berg, S. Patil, and R. Alterovitz, "Efficient approximate value iteration for continuous gaussian POMDPs." in AAAI, 2012.
- [13] D. Wierstra, A. Foerster, J. Peters, and J. Schmidhuber, "Solving deep memory POMDPs with recurrent policy gradients," in *ICANN*, 2007.
- [14] S. Omidshafiei, C. Amato, M. Liu, M. Everett, J. P. How, and J. Vian, "Scalable accelerated decentralized multi-robot policy search in continuous observation spaces," in *ICRA*, 2017.
- [15] R. Sabbadin, J. Lang, and N. Ravoanjanahry, "Purely epistemic markov decision processes," in AAAI, 2007.
- [16] C. Guestrin, D. Koller, R. Parr, and S. Venkataraman, "Efficient Solution Algorithms for Factored MDPs," *JAIR*, 2003.
- [17] H. Kurniawati, D. Hsu, and W. Lee, "SARSOP: Efficient Point-Based POMDP Planning by Approximating Optimally Reachable Belief Spaces," in RSS, 2008.
- [18] G. R. Bradski, "Computer vision face tracking for use in a perceptual user interface," in *WACV*, 1998.
- [19] Y. Watanabe, A. Manecy, A. Amiez, C. Lesire, and C. Grand, "Noncooperative ground vehicule tracking and interception by multi-rpa collaboration," in *ICAS*, 2016.
- [20] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," ACM T-IST, vol. 2, no. 3, 2011.
- [21] J. Pineau, G. Gordon, and S. Thrun, "Point-based value iteration: An anytime algorithm for POMDPs," in *IJCAI*, 2003.
- [22] N. Drougard, F. Teichteil-Konigsbuch, J.-L. Farges, and D. Dubois, "Qualitative possibilistic mixed-observable MDPs," in UAI, 2013.