

Active Object Detection and Pose Estimation in General Belief Space

Kanzhi Wu, Teng Zhang, Gamini Dissanayake, Shoudong Huang and Ravindra Ranasinghe
Centre for Autonomous Systems, University of Technology Sydney

kanzhi.wu, teng.zhang, gamini.dissanayake, shoudong.huang, ravindra.ranasinghe@uts.edu.au

I. INTRODUCTION AND RELATED WORK

3D object detection and pose estimation is an indispensable task for service robots, however, due to the challenges in everyday environments where occlusion are observed frequently, it has been recognised that object detection and localisation cannot be achieved using observation from a single viewpoint. Active object detection, recognising the targeting objects via active maneuver, is a solution towards improving the recognition probability and minimising the estimation uncertainty through fully utilising the mobility of the robotic platform. Active perception, proposed by Bajcsy[1], is defined as the problem of intelligent control strategies applied to the data acquisition process which will depend on the current state of data interpretation including recognition. Active object detection and pose estimation, as a subset of active perception, has been investigated intensively in recently years. Atanasov et al. [3] introduced Viewpoint-Pose tree (VP-Tree) which provides not only object category but also its relative pose and the active hypothesis testing problem is addressed by a point-based POMDP algorithm. Potthast et al. [2] combined next-best-view planning and feature selection into one online active object detection system. Based on simple features and detectors, their algorithm enables the quadrotor to detect the targeting object more robustly.

In this paper, discretisation on state and control spaces are not required and similar as Indelman's work[4], based on the probabilistic representation of the general belief, the next-best-view are further decided via optimising the objective function parameterised by the general belief, thus enabling the planning in the continuous domain. We follow the *maximum likelihood observation* assumption and construct the objective function which consolidates uncertainties of the robot and the poses of detected objects, object detection confidences and control consumption. The simulation results demonstrates the effectiveness of the proposed approach.

II. NEXT-BEST-VIEW DECISION IN GENERAL BELIEF SPACE

A. Preliminaries and problem formulation

The parameters which will be used in following sections are listed in Table I and the noises appeared in the table are assumed to be Gaussian. In Table I, the general motion model and observation model can be further detailed. For motion model G , without acquiring observation, object states

TABLE I: Summary of Notations

Symbol	Description
X^o	object pose
X^r	robot pose
\mathbf{X}_k	$\mathbf{X}_k = [X_1^o, \dots, X_i^o, X_1^r, \dots, X_k^r]^\top$ which is the state of time step k , i is the number of detected objects up to time step k
$z_{k,j}$	observation of feature j at time step k . Multiple features are distributed on each object and its coordinate w.r.t the object coordinate frame is given
Z_k	$Z_k = [z_{k,1}, z_{k,2}, \dots, z_{k,n_k}]$ where n_k is the number of observed features at time step k
$Z_{1:k}$	$\mathbf{Z}_k = [Z_1, Z_2, \dots, Z_k]^\top$
u_{k-1}	control input at time step $k-1$
$u_{1:k-1}$	$\mathbf{u}_{k-1} = [u_1, u_2, \dots, u_{k-1}]$
\mathbf{G}	$X_k = \mathbf{G}(X_{k-1}, u_{k-1}) + \eta$ is the motion model
\mathbf{H}	$Z_k = \mathbf{H}(X_k) + \xi$ is the observation model

$[X_1^o, X_2^o, \dots, X_i^o]$ in \mathbf{X}_k will not be updated and only the robot pose X^r will be computed according to:

$$X_k^r = g(X_{k-1}^r, u_{k-1}) + \eta_r \quad (1)$$

Also, the observation Z_k is only decided by the current pose of the robot X_k^r and estimated poses of the objects $[X_1^o, \dots, X_i^o]$ thus \mathbf{H} can be simplified as:

$$Z_k = \mathbf{H}(X_1^o, \dots, X_i^o, X_k^r) + \xi \quad (2)$$

Considering each feature individually, we have:

$$z_{k,j} = h(X^o, X_k^r, p_j^f|_o) + \xi_j \quad (3)$$

where $p_j^f|_o$ is the coordinate of feature j in object frame. Similar as Indelman's work, the generalised belief in l -th predication step is denoted as:

$$\mathbf{b}(X_{k+l}) \doteq p(X_{k+l} | Z_{1:k}, u_{1:k-1}, Z_{k+1:k+l}, u_{k:k+l-1}) \quad (4)$$

and this belief is assumed to be Gaussian as:

$$\mathbf{b}(X_{k+l}) \sim \mathcal{N}(X_{k+l}^\Delta, I_{k+l}^{-1}) \quad (5)$$

where X_{k+l}^Δ is the *maximum a posterior* (MAP) estimation and I_{k+l} is the corresponding information matrix. In order to find the optimal control inputs $u_{k:k+L-1}^\Delta$ in next $L-1$ steps, we need to minimise the objective function J_k parameterised by $u_{k:k+L-1}$ and the detailed formulation of J_k will be illustrated in the following section.

$$u_{k:k+L-1}^\Delta = \arg \min_{u_{k:k+L-1}} J_k(u_{k:k+L-1}) \quad (6)$$

B. Objective function formulation

In order to compute the optimal control $u_{k:k+L-1}^\Delta$ via optimising the objective function J_k , we can divide the problem into 2 steps: 1) inference $\mathbf{b}(X_{k+L})$ given control inputs $u_{k:k+L-1}$ and 2) compute the $u_{k:k+L-1}^\Delta$ which minimises objective function. The 2 steps are executed iteratively until the stop criterion are met. The second step can be realised by gradient based optimisation methods and the control inputs will be updated as:

$$u_{k:k+L-1}^{(i+1)} = u_{k:k+L-1}^{(i)} - \lambda \nabla J_k \quad (7)$$

where ∇J_k is the gradient of objective function at guess $u_{k:k+L-1}^{(i)}$. The main difficulty lies in predicting the general belief $\mathbf{b}(X_{k+l})$ without acquiring the observations. Following by the derivations in [4] but with slightly difference (assuming maximum likelihood observation), using *Expectation-Maximisation* and formulating the objective function as a quadratic form, we approximate the X_{k+l}^Δ and $I_{k+l}(u_{k:k+l-1})$:

$$\begin{aligned} X_{k+l}^\Delta &= \bar{X}_{k+l} \\ I_{k+l} &\doteq \mathcal{A}_{k+l}^\top \mathcal{A}_{k+l} \quad \text{where } \mathcal{A}_{k+l} = \begin{bmatrix} I_k^{1/2} & \mathbf{0} \\ \mathcal{G} \\ \tilde{\mathcal{H}} \end{bmatrix} \end{aligned} \quad (8)$$

where \bar{X}_{k+l} is the nominal state and \mathcal{G} and $\tilde{\mathcal{H}}$ are stacked by the Jacobians for motion model and observation model.

The objective function shown below consists 3 terms (costs):

$$\begin{aligned} J_k(u_{k:k+L-1}) &= \text{Tr}(M_\Sigma I_{k+L}^{-1} M_\Sigma^\top) + \sum_{i=1}^{n_{obj}} \Psi(\Upsilon|_{n_i \times (k+L)}) + \\ &\quad \sum_{l=0}^{L-1} \Phi(u_l) \end{aligned} \quad (9)$$

- 1) Uncertainties of the robot pose and the poses of detected objects. In eq. 9, the uncertainty is obtained by calculating the trace(weighted trace) of the covariance matrix as $\text{Tr}(M_\Sigma I_{k+L}^{-1} M_\Sigma^\top)$;
- 2) Object detection confidences parameterised by observed features. Based on the *MAP* estimate of X_{k+L} , we construct a *observation score matrix* $\Upsilon|_{n_i \times (k+L)}$ which describes the feature matching score under a specific relative pose where n_i is the number of features in object i . We compute the *observation score* $v = \Psi(\Upsilon)$ which describes the probability of missing the targeting objects after a sequence of observations;
- 3) Control constraints where larger control inputs will be penalised. In eq. 9, the cost function is denoted as $\Phi(u_l)$. In Section. III, we assign larger cost for larger angular velocity control input;

III. EXPERIMENTS AND RESULTS

We present simulation results, the planned trajectories for active object detection, using the proposed approach in fig. 1. The first row in fig. 1 shows the results for single object and the

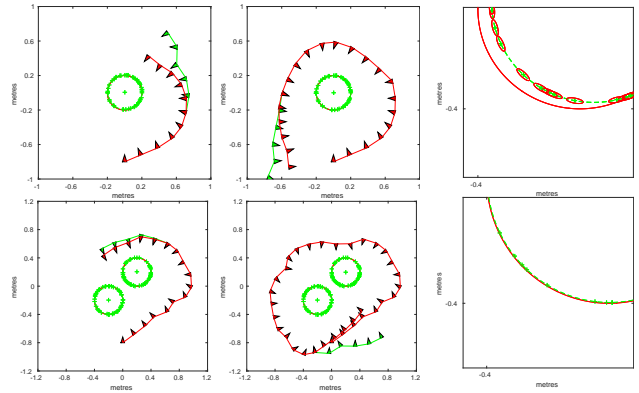


Fig. 1: Planned trajectories for active object detection

Fig. 2: Uncertainty of targeting object. The second row shows the planned trajectories for multiple objects scenario. The green trajectories are computed from the initial control input and red ones are generated from the optimised control input and executed. During the simulation experiments, we consider the next 5 steps in the planning phase while only one step is executed. Therefore, for each term in the objective function, larger weight is associated with nearer steps which indicates that the robot are more interested in short-term gain.

The state, \mathbf{X}_k , is updated using EKF after acquiring the updated observation and the uncertainty are reduced after observing the object for multiple times. Fig. 2 demonstrates how the uncertainty of targeting object are reduced after multiple observations during active object detection. In the upper figure in fig. 2 which shows the initial uncertainty of the object, ellipses represents the uncertainty of the features on the object, the dashed green circle denotes the current estimation of the object pose and the red circle denotes the ground-truth pose. The bottom figure in fig. 2 shows the converged object pose and features where the sizes of the red ellipses are significantly reduced and the estimated object pose is closer to the ground-truth.

This paper presents a probabilistic framework for active object detection and pose estimation which is able to provide the optimal control in continuous space. The remaining difficulty is to eliminate the maximum likelihood observation assumption thus considers the unexpected occlusion or noise during the prediction. The experiments on actual robotic platform will be conducted in future work.

REFERENCES

- [1] R. Bajcsy. Active perception. *Proc. of the IEEE*, 76(8): 966–1005, Aug 1988. ISSN 0018-9219.
- [2] Christian Pothast et. al. Active multi-view object recognition and online feature selection. In *ISRR*, Sep 2015.
- [3] N. Atanasov et. al. Nonmyopic view planning for active object classification and pose estimation. *TRO*, 30(5): 1078–1090, Oct 2014.
- [4] Vadim et. al Indelman. Planning in the continuous domain: A generalized belief space approach for autonomous navigation in unknown environments. *IJRR*, 34(7):849–882, 2015.