

Active Multi-View Object Recognition and Change Detection

Christian Potthast, Andreas Breitenmoser, Fei Sha and Gaurav S. Sukhatme

I. INTRODUCTION

In many systems, especially in robotic systems, active reasoning and decision making is essential for a successful deployment. The goal of active reasoning is to choose subsequent actions with the highest utility to maximally increase the reward gained by executing the new actions. As a result, this can improve system performance drastically while at the same time decreasing computation time. In general, every multi-action system can benefit this way, but the benefits become more apparent when the system has minimal computational power and only a limited amount of operation time available. This is typically the case for resource-constrained mobile robots such as quadcopters.

With regard to single-view object recognition methods, many have been developed over the years, but despite steady improvements, no perfect system has yet been found. Many of the developed systems have complex and computationally expensive feature representations, making them difficult to use on a robotic system with limited computational power. Feature representation is in particular hard, because we need suitable features that generalize well across different object categories but remain expressive enough for accurate object classification. Moreover, especially single-view object recognition is affected by object ambiguity, i.e., two objects which look very similar are hard to tell apart.

To overcome those difficulties, we propose a *multi-view* Bayesian framework that performs *active view planning* [1] and *online feature selection* [2]. Furthermore, we show how this multi-view recognition system can be used for object change detection. Given a prior map of the environment, the task is to determine whether objects in the environment have changed. This can either mean that the object has changed completely (i.e., it belongs to a different object class), or the object has changed its orientation only (i.e., it belongs to the same class but its rotation differs in yaw).

We extensively evaluate our active reasoning system for the two perception tasks of *object recognition* and *object change detection* on a large RGB-D dataset [3] and show first preliminary results from deploying the system on a quadcopter robot.

II. MULTI-VIEW OBJECT RECOGNITION

Our multi-view object recognition framework consists of a Bayesian network and a sequential update procedure that allows to integrate new observations and to infer about the

University of Southern California (USC), Department of Computer Science, Los Angeles, CA 90089, USA {potthast, andreas.breitenmoser, feisha, gaurav}@usc.edu



Fig. 1. Test setup for active multi-view object recognition and change detection. RESL quadcopter equipped with an ASUS Xtion RGB-D sensor, flying around a target object.

new posterior distribution over all object classes. After each integration step at time t , we take another action a_{t+1} and compute a new feature f_{t+1} . This can either be 1) extracting a new observation feature at the current location and updating the posterior distribution, or 2) moving the sensor to a new location, where a new feature is evaluated. To find the next best possible action, we compute a utility score s_t and choose the next action and feature as $\{a_{t+1}, f_{t+1}\} = \arg \max_{a, f} s_t(a, f)$, given all the possible actions a and available features f in the feature set.

In an offline training phase, we first train a set of generative object models $o^k = \{c^k, \theta^k\}$, with $k \in \{1, \dots, K\}$, c^k being the object class and θ^k the object orientation angle. These models are trained from N different feature types $\mathbf{f} = [f^1, \dots, f^N]$; each resulting feature vector is represented in the model as a combination of independent Gaussian distributions, $p(f|o^k) = \mathcal{N}(f|\mu_{o^k}, \sigma_{o^k})$. Given the models, we can infer about the observed features $f_{1:T}$ and taken actions $a_{2:T}$ at time $t = T$, and compute the object recognition posterior distribution $P(o_T|a_{1:T}, f_{1:T})$.

The utility score for a particular action can be computed by the *mutual information* (MI), which computes the reduction in uncertainty of the current object class and pose if a new observation is made. We use MI for computing the utility score of an action a to move to a new viewpoint by

$$s_t(a, f) = I(o_t; f|a) \approx \sum_{k=1}^K P(o^k|f_{1:t}) \cdot P(f_{\mu_{o^k}}|o^k, a) \cdot \log \frac{P(f_{\mu_{o^k}}|o^k, a)}{\sum_{i=1}^K P(f_{\mu_{o^i}}|o^i, a)}.$$

Usually, we would have to marginalize over the feature distributions because we do not know exactly what the features look like. However, for efficiency we compute MI with a MAP approximation, by sampling with zero variance from the feature distributions, resulting in the overall mean feature vector \mathbf{f}_μ .

Method	Accuracy	Observations
C+B+S+V (no FS)	84.48% \pm 1.4	3.1 \pm 0.2
Random	86.76% \pm 0.4	3.2 \pm 0.2
Mutual Information	91.87% \pm 0.7	2.5 \pm 0.1

Table I: Active multi-view object recognition. The first row uses the full feature set with random viewpoint selection; the next two rows use feature selection and the respective view planning method: random or MI.

Similar to computing the utility score for a new viewpoint above, we can also use MI to compute the utility of adding a new feature f from the set of all possible features. This allows us to only select and infer about relevant features, marginalizing out unnecessary features without loss of information.

III. CHANGE DETECTION

We now show how to utilize our framework for detecting changes in the environment. Given a prior map of the environment and an object database, we want to detect if objects in the environment have been replaced, taken away or changed their pose. We use our generative object model to express the probability $P(o^k|f)$ that the object has been generated by this feature. Given an expected object model \hat{o}^k (prior information), we can compute the probability $P(u)$ that the object has changed with $P(u) = 1 - P(\hat{o}^k|f)$. If the probability $P(u) > \tau$, we know that the feature observed cannot be generated by the model \hat{o}^k , hence either object class or object pose has changed. This can also occur when the feature does not match the model anymore due to changes in lighting or noise; in either case we need to perform further observations to definitively conclude about the change. If we are uncertain about the object after our first observation, we compute new viewpoints that allow for observations that are unique in terms of feature space. We acquire new observations, evaluate additional features and compute new likelihoods of the features given the expected object model \hat{o}^k .

IV. RESULTS

We evaluate our active multi-view object recognition and change detection framework on the large dataset described in [3]. The dataset consists of 300 objects, captured from all sides and three different observation positions on the vertical axis with viewing angles of 30°, 45° and 60°. We train our object models using two of the observation positions (30°, 60°) and test on the data captured at 45°. During test time, we want to predict the correct instance of the object class by taking the MAP estimate of the posterior distribution.

The object models are trained from the feature vector f , which contains four independent features: object bounding box, color, SIFT and the geometric feature VFH. Although these features are all fairly simple, they still result in expressive recognition results in a multi-view setting as we can see in Tab. I. The table shows the performance of our multi-view object recognition framework in three different settings. The first row represents the base case, new viewpoints are chosen at random and the full set of available features is used, meaning no feature selection is applied. The next two rows evaluate the performance of feature selection in combination

Method	Accuracy	Observations
No change	98.3% \pm 0.21	1.2 \pm 0.19
Pose change	99.4% \pm 0.12	1.3 \pm 0.22
Object change	92.4% \pm 0.15	1.9 \pm 0.27

Table II: Active multi-view change detection with quadcopter robot. The accuracies are shown for detecting no change in an object (first row), change in pose only (second row) and complete object change (third row).

with two different viewpoint selection methods, *random* selection or MI. From row 1 and 2 we can see that feature selection increases the accuracy slightly, however, the main advantage of feature selection lies in the reduction of computation time. By only computing and integrating features that add information, we save valuable time. Unfortunately, due to limited space, we have to omit here these further results. Finally, we can see a big performance increase by using feature selection in combination with intelligent viewpoint selection, achieving a respectable 91.87% in recognition accuracy and needing less observations on average.

Moreover, we present first preliminary results of our object change detection system on our quadcopter platform (Fig. 1). The quadcopter, equipped with an RGB-D sensor, is tasked to identify the change of a target object by autonomously selecting observation positions and infer about the captured data. Since we do not own the original objects of the large dataset, we added ten objects similar to the ones already in the dataset, resulting in 310 total objects, and tested on the new objects. Three experiments have been performed: 1) the target object does not change, 2) the object’s pose changes and 3) the object changes, with the results averaged over eight different starting locations.

In Tab. II we show the results of detecting whether change has occurred to an object. Detecting that there is no change is the hardest case, since it is very similar to recognizing the object class and pose. Detecting changes in pose or that the object has changed completely is easier because we test a feature that is very unlikely to be generated by the expected model. On average, we can see that in all three methods we need less than two observations. However, in some instances we are uncertain after only one observation and we need to take an additional one.

V. CONCLUSION

We have presented an active multi-view object recognition and change detection framework which incorporates an information-theoretic approach for active viewpoint selection and feature selection. Our experiments have shown that we can achieve respectable results by using a multi-view approach in combination with relatively simplistic features. Furthermore, we have demonstrated our multi-view approach on a quadcopter robot to detect object change.

REFERENCES

- [1] N. Atanasov, B. Sankaran, J. L. Ny, G. J. Pappas, and K. Daniilidis, “Nonmyopic View Planning for Active Object Classification and Pose Estimation,” *IEEE Transactions on Robotics*, vol. 30, no. 5, 2014.
- [2] M. Verleysen and F. Rossi, “Advances in Feature Selection with Mutual Information,” *CoRR*, 2009.
- [3] K. Lai, L. Bo, X. Ren, and D. Fox, “A large-scale hierarchical multi-view RGB-D object dataset,” in *ICRA*. Shanghai: IEEE, 2011.