Next best view estimation with eye in hand camera

Christian Potthast and Gaurav S. Sukhatme

Abstract—

I. INTRODUCTION

In this work we are interested in automatic data acquisition using a high dimensional robotics arm equipped with a camera system. We want to iteratively estimate new view points that maximizes a reward function and simultaneously optimize the trajectory to reach the new view point given a second reward function. For many robotic systems automated data acquisition is important to detect objects, recognize objects or simply estimate object pose. The Next-Best-View (NBV) problem, first discussed by [1], seeks a single additional sensor placement in order to improve the knowledge of the environment. This can be formulated as a local optimization problem with constraints. Depending on the context the optimization problem has different constraints. Cluttered environments impose visibility constraints due to occlusion thereby increasing the difficulty of making good predictions of the next best view location. Optimize a trajectory from a start to a given goal pose is a high dimensional planning problem. Because of the high dimensionality we have to sample the space for suitable waypoints. The challenge is to sufficiently sample the space to find good waypoint and therefore yield a good trajectory.

The increase in degrees of freedom allows the robot to choose from a bigger set of possible viewing positions compared to robots which can only move along on a 2D plane. Oftentimes obstacles, like walls, furnitures, etc., block viewing positions which prevents the robot of gathering more information about a certain area of interest. The extension of possible viewing positions due to an increase in degrees of freedom raise the changes of observing an area despite view blocking obstacles. A highly maneuverable arm can advance in areas which otherwise would have been left unobserved.

Previous developed techniques that used eye in hand cameras [2] often times neglect the fact that while the manipulator is positioning itself to a new position data can be gathered. This lead to the effect that the problem is strictly handled as a high dimensional next best view planning problem. Additionally most of the prior work [1] and [3] is not dealing with cluttered environments. Such problems are oftentimes solved using a geometric approach, computing sensor positions to view a desired space. This approach will fail in cluttered environments since hidden objects in unknown space can potentially block the desired space.



Fig. 1. PR2 robot equipped with a Microsoft Kinect in the right manipulator.

In our approach to the cluttered environment we reason about the unknown space using probabilistic methods. As a result, our algorithm chooses the next best scanning position that promises the highest expected information gain. Additionally we are using a Probabilistic Road Map planner [4] to plan a path from a given start position to a goal position in a high dimensional space. The path with the highest reward is execute to move the manipulator to the next best view position.

Our experimental platform is the PR2 from Willow Garage shown in Fig. 1. The PR2 robot is a research and development platform which uses the open source robot operating system (ROS). The PR2 is equipped with two 7 DOFs manipulators. For data acquisition we place one Microsoft Kinect sensor into one of the robots manipulators.

II. RELATED WORK

In general, the NBV problem can also be stated as a view planning problem. The objective is to find a set of sensor configurations which fulfill a modeling or surveillance objective. Automated 3D scene planning and data acquisition employing robotic systems are discussed in [5] and [6] for indoor environments. A mobile robot for outdoor scene modeling of large structures (using a voxel grid approach in combination with a ray traversal algorithm to estimate potential next positions) is discussed in [7]. Prior knowledge is used in form of a 2D map which specifies the footprint of walls. This footprint is used to sample new positions for estimating the next best position by counting voxels which are on the boundary from unknown to free. In [3] a probabilistic surface reconstruction approach is described

The Authors are with the Robotic Embedded Systems Laboratory (RESL) and Department of Computer Science at the University of Southern California (USC), Los Angeles, CA 90089, USA {potthast, gaurav}@usc.edu

using depth cameras. Given a current surface reconstruction a metric is formulated to score potential new scanning positions in terms of reduction of uncertainty from virtual created depth maps. A eye in hand camera is used in [2] to model an object in the environment by searching in a high dimensional space for the next best view position.

In contrast to previous work, we reason over the unknown space itself and estimate the expected information gain in cluttered environments. Additionally we incorporate measured data while moving the sensor to the next best view position. Our work is closest to [1], [3] and [2]. However, all three [1], [3] and and [2] are not well suited for cluttered environments since they assume full visibility to the target region or object. Unseen objects in the occluded space lead to overestimation of the reward for potential future scanning positions. Furthermore all approach do not take into account the information gathered while the sensor is moving to the next view position.

III. TRAJECTORY PLANNING

We are given a partially observed scene in an area of interest, which we would like to further investigate using a camera mounted on a robot manipulator. We want to find view points which globally minimize the number of unknown space in the environment and find trajectories which locally increase knowledge about objects. To estimate the next best view position as well as a good trajectory we are using a PRM approach. In the two phases of the PRM, which are the planning and query phase. In the planning phase we estimate a globally next best view and build the graph while in the planning phase we plan a path which locally optimizes the increase in object knowledge.

A. PRM planning phase

In the planning phase of the PRM we want to estimate the next best viewing position. This is done by randomly sample new poses, evaluate the potential decrease in unknown space and connect the newly sampled poses with the k-nearest neighbors to build a connected graph. We use a sampling approach because the cartesian workspace of the manipulator is high dimensional and we have no precise knowledge about the environment we can not simply use a geometric approach and compute the next best viewing pose.

To create the PRM we randomly sample a 6D pose $s \in \mathbb{R}^3$, with $S = \{s_1, s_2, \ldots, s_n\}$ in the cartesian manipulator workspace. Before inserting the newly created pose as a node into the graph we have to check its feasibility in terms of reachability and joint constraints using inverse kinematics. After the feasibility check each new pose is inserted into the graph by connecting the new node to its k-nearest neighbors. The distance metric to determine the nearest neighbors is the euclidean distance. However in general one would rather prefer a distance metric in joint space to prefer joint space distances.

In order to reduce unknown space we need to minimize the number of unknown voxels. From this point of view observing an occupied voxel is equally valuable as detecting



Fig. 2. A PRM created in the cartesian manipulator workspace with the red arrows being the random sampled poses and the. The green arrow is the computed the global next best view pose. The best path which optimizes the local information gain is drawn as a green line connecting the start pose with the end pose.

a free voxel. Voxel state changes are therefore a measure of information gain in a new scan. Thus, minimizing the unknown space is equivalent to maximizing the expected number of observed voxels. Intuitively, we aim to find a new viewing position which has the largest information gain compared to all other possible positions. To find this new position we have to make a guess about what we can observe from all possible positions and predict how the voxel grid would change if we were to take a scan from each position. This is achieved by taking 'virtual scans' from every possible position and then choosing that position that promises the most information gain.

The next best view pose \hat{s} is estimated by choosing the pose s_n with the maximum expected number of observed voxels, where the function $F(s_n)$ is a measure of how many unknown voxels can be observed from s_n .

$$\hat{s} = \operatorname*{arg\,max}_{n} \mathbb{E}\big[F(s_n)\big] \tag{1}$$

The exact method of how to estimate the next best viewing pose is describe in [8].

B. PRM query phase

In the query phase of the PRM we want to estimate a path given a start pose P_s and a goal pose P_g . The start pose P_s is the current robot manipulator pose and the goal pose P_s is the estimated next best view pose \hat{s} hence $P_g = \hat{s}$. We want to find a trajectory which maximizes the expected number of observed voxels which are close to objects. These voxels are voxels which belong to objects, by observing them our knowledge about the object itself will increase.

To find the best path in our graph we have to estimate all paths $r = \{r_1, r_2, \ldots, r_n\}$ which connect the start pose P_s and a goal pose P_g . To find the best path in the graph we have to evaluate each node in the graph of their local information gain. The expected local information gain is computed by estimating

$$l = \mathbb{E}\big[L(s_n)\big] \tag{2}$$

where the function $L(s_n)$ is a measure of how many unknown voxels close to objects can be observed from s_n . The best path \hat{r} is computed by choosing the path r which maximizes the expected local information gain l in the sequence of nodes in r. The best path is then computed with

$$\hat{l} = \operatorname*{arg\,max}_{n} \sum_{n} \mathbb{E} \big[L(s_n) \big] \tag{3}$$

The information gain l of the best path \hat{r} is either increasing or keeps constant if no new information is gain by moving the manipulator to a new pose. The nodes of the best path \hat{r} are then taken as waypoints for the manipulator trajectory. A constructed PRM with the best path computed can be seen in Fig. 2.

IV. PREPROCESSING

The general framework of the preprocessing steps can be found in [8]. We however have used a occupancy grid representation of the environment compared to our previous work. In our previous work we have assumed a perfect measurement model which we have replaced in this work with a gaussian noise model. The occupancy grid is cable to probabilistic model the environment given such a noise model.

A. Occupancy grid

To estimate and represent unknown space in the environment we use a occupancy grid in combination with a Bayesian voxel update formulation, which can be seen in Fig. 3. Every voxel $v_{x,y,z}$ in the occupancy grid \mathcal{O} is assigned a occupancy probability P(v). Initially, all voxels $v_{x,y,z}$ of \mathcal{O} are initialized with P(v) = 0.5, to state that we have no knowledge yet about the scene. Our measurements are given as a point cloud \mathcal{P} . For every new measurement $p \in \mathcal{P}$, we want to update the corresponding voxel $v_{x,y,z} \in \mathcal{O}$ with a new occupancy probability.

To do so we are using a stochastic sensor model defined by a probability density function p(m|z), with m being the actual measurement and z being the true parameter state space value. In our case z corresponds to the cartesian coordinate of the voxel center. The probability function p(m|z) is gaussian

$$\mathcal{N} = (\mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-1/2(r-d)^T \Sigma^{-1}(r-d)\right)$$

To determine the cell occupancy probability we want to incrementally update the sensor information into the occupancy grid using sequential update formulation of the Bayes' theorem. Given a current estimation of the occupancy probability $v_{x,y,z}$, we want to estimate $p(v_{x,y,z}|m_{x,y,z})$, given a new measurement m. The update is formulated with

$$p(v = OCC|m_{t+1}) = \frac{p(m_{t+1}|p(v) = OCC) * p(p(v) = OCC|m_t)}{\sum_{p(v)} p(m_{t+1}|z) * p(z|m_t)}$$



Fig. 3. The environment represented as an occupancy grid map.

for each voxel in the occupancy grid. Although we want to possibly keep the probabilistic representation of the environment for now we threshold the occupancy probability and assign every voxel one of three states:

$$v_{x,y,z} = \begin{cases} 1, & \text{free space} \\ 0, & \text{occupied space} \\ -1, & \text{unknown space} \end{cases}$$
(4)

For completeness, let the state s be assigned after following thresholding:

$$s = \begin{cases} p(v) <= 0.4, & 1\\ p(v) >= 0.6, & 0\\ p(v) < 0.6 \land p(v) > 0.4, & -1 \end{cases}$$
(5)

1) Unknown space estimation: The point cloud produced from the Kinect sensor gives us no information about the space where no measurements have been returned. The point cloud itself does not distinguish between space which is free and space which could not be measured due to occlusion or range limitations. We need to distinguish between the two types of space, because the unknown space arising from occlusion is the space we are interested in. This unknown space can give us new information about the scene once it is observed. By using a ray traversal algorithm we can distinguish between space we have seen and is actually free and space which we have not seen and is therefore unknown space.

To estimate the occluded part of the scene we shoot rays from position P to every uncertain voxel $u_{x,y,z} \in \mathcal{U}$ in the voxel grid. P is the position of the sensor from which it has acquired the input point cloud. The state of the voxel depends on the ray ending in the voxel $u_{x,y,z} \in \mathcal{U}$. Every ray rcontains all voxels it penetrates until it hits the target voxel. The voxels which are penetrated by the ray influences the state of the target voxel. If the ray contains a voxel $v_{x,y,z} = 0$ (meaning the ray has penetrated a voxel which is occupied) the state of the target voxel will remain at state -1. If the ray contains only free or unknown voxels we can change the state of the voxel to 1 as being free. This state change is



Fig. 4. The PR2 is estimating the next best viewing pose and simultaneously estimate a path with locally optimizes the increase in object information gain. The PR2 is then moving the manipulator equipped with the Kinect camera to the new goal pose following the best estimated path.

justified because P_i was the actual scan position and we can be certain that we see the unknown voxel if we penetrate only free voxels or unknown voxels. Otherwise we would have measured an occupied voxel and therefore penetrated an occupied voxel.

V. CONCLUSIONS

We have implemented a fully automated system on the PR2 to estimate the next best viewing pose for a eye in hand camera. To choose the next viewing pose we develop an information gain-based variant of the next best view problem. Furthermore a path is computed to the next best view pose which maximizes the local information gain to increase object knowledge. To estimate the best trajectory we have used a PRM approach to sample the cartesian workspace of the robot and find the trajectory that maximizes the local information gain.

VI. FUTURE WORK

Our system has still certain limitation and a lot of room for improvements. I will list some of them as future work.

A. Adaptive Planning

An extension to the previous described method is to adaptively change the trajectory if predictions and real world scans while the manipulator is moving towards the goal are differ. Predicting the outcome for every node in the tree has the effect that errors in the calculated expected number of seen voxels yield to either over or under estimation with an increase in lookahead. This has the effect that by executing the computed 'best' trajectory, errors in the prediction, lead to non 'best' trajectories. Because space could have been already earlier seen as expected or objects appear which have not been expected. To prevent this we propose to adaptively change the trajectory if the error between predicted number of seen voxel and the real captured scan is bigger as a certain threshold.

B. Occupancy grid

Right now the occupancy grid is thresholded to estimate their voxel states. This has the effect that the probabilist representation of the environment is lost. A further improvement would be to not threshold but rather use a probabilistic ray traversal algorithm which assign probabilities to unknown voxels instead of a state.

C. Next Best View Estimation

The next best view estimation is very dependent on the clutter of the environment. A more general approach could improve the calculation of the expected number of visible approach. One could also think about an adaptive approach which changes after some measurements have been taken. This however requires some kind of a measure for clutter.

REFERENCES

- C. Connolly, "The determination of next best views," *Robotics and Automation. Proceedings.* 1985 IEEE International Conference on, vol. 2, pp. 432 435, 1985.
- [2] K. G. Liila Torabi, "Integrated view and path planning for an autonomous six-dof eye-in-hand object modeling system." IEEE/RSJ International Conference on Intelligent Robots and Systems, 2010.
- [3] M. Krainin, B. Curless, and D. Fox, "Autonomous generation of complete 3d object models using next best view manipulation planning," in *IEEE international conference on Robotics and Automation (ICRA)*, 2011.
- [4] L. E. Kavraki, P. Svestka, J.-C. Latombe, and M. Overmars, "Probabilistic roadmaps for path planning in high dimensional configuration spaces," *IEEE Transactions on Robotics and Automation*, vol. 12, no. 4, pp. 566–580, 1996.
- [5] A. Nüchter, H. Surmann, and J. Hertzberg, "Planning robot motion for 3d digitalization of indoor environments," *ICAR*, Jan 2003.
- [6] Z. R. S. J. Eidenberger, R., "Probabilistic occlusion estimation in cluttered environments for active perception planning," in Advanced Intelligent Mechatronics, 2009. AIM 2009. IEEE/ASME International Conference on. Advanced Intelligent Mechatronics, 2009.
- [7] P. Blaer and P. Allen, "Data acquisition and view planning for 3-d modeling tasks," *IEEE/RSJ IROS*, Jan 2007.
- [8] C. Potthast and G. Sukhatme, "A probabilistic framework for next best view estimation in a cluttered environment." IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2011, (Submitted).