

Segmentation and Classification of Range Image from an Intelligent Vehicle in Urban Environment

Xiaolong Zhu, Huijing Zhao, Yiming Liu, Yipu Zhao, Hongbin Zha

Abstract—As the rapid development of sensing and mapping techniques, it becomes a well-known technology that a map of complex environment can be generated using a robot carrying sensors. However, most of the existing researches represent environments directly using the integration of point clouds or other low-level geometric primitives. It remains an open problem to automatically convert these low-level map representations to semantic descriptions in order to effectively support high-level decision of a robot. Based on another representation of 3D point clouds, i.e. range image, this paper proposes a framework of segmentation and classification of range image, the objective of which is to annotate class labels to the data clusters that are obtained through a graph-based segmentation. Experimental results are presented and evaluated demonstrating that the proposed algorithm has efficiency in understanding the semantic knowledge of a large dynamic urban outdoor environment.

I. INTRODUCTION

As the rapid development of localization and mapping techniques, e.g. SLAM (Simultaneous Localization And Mapping), it becomes a well-known technology that a map of a complex environment can be obtained by a robot carrying sensors [21]. Most of the maps, as many researchers successfully demonstrated [6], [8], [18], [23], are represented as an integrations of 3D points, or other low-level primitives, such as feature points, lines, planar surfaces and so on. However, such a representation tells only existence of spatial entities. A robot can not directly retrieve from the data what kind of objects are there in the surroundings. In order to support more intelligent decision-making, an automatic tool is needed for a robot to convert those low-level map representations to higher-level ones that contain semantic knowledge of the scene. The contribution of this work is to propose a system, which formats the 3D data into a range image so that standard vision techniques could apply to segmentation and further classification.

In our previous research, a vehicle robot [26], called POSS-v, has been developed. As shown in Fig.1, five single-row laser scanners are mounted on the vehicle profiling the surroundings from different viewpoints at different directions. Laser scanner L1 conducts horizontal scanning, the data of which is mainly used, along with the GPS/IMU data, to perform a SLAM, so that vehicle pose with both local and global accuracy can be achieved in dynamic urban environment [27]. Other four laser scanners, L2-L5, are used to acquire map data of road surface, objects above the road,

and objects to the right or left of the vehicle respectively. As the vehicle moves along streets, a range image is obtained from each scanner, which can be converted to a 3D point cloud in a global coordinate system by integrating the parameters of both sensor geometry and vehicle pose.



Fig. 1. An intelligent vehicle (POSS-v) with multiple single-row laser scanners.

This paper focuses on processing of the data from laser scanners L4 and L5, which are vertical profiling to the objects on road sides, using range images as the data representation. A clip of range image is shown in Fig.2, where each column corresponds to a laser scan line, and each pixel is represented by converting the laser range to an intensity value in [0,255]. For each range value, it has a 2D coordinate in range image, and also a 3D coordinate in a global coordinate system. All the information is used in the following data processing.

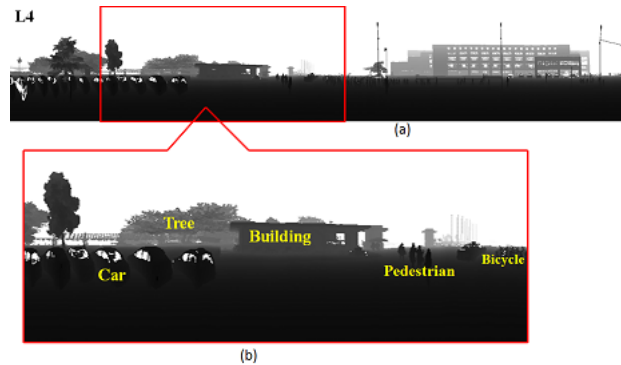


Fig. 2. (a) a clip of range image obtained by L4, (b) enlarged figure with some of the objects annotated manually.

This work is partially supported by the NSFC grants (No. 90920304 and No. 60975061).

X.Zhu, H.Zhao, Y.Liu, Y.Zhao and H.Zha are with the State Key Lab of Machine Perception (MOE), Peking University. {zhaohj}@cis.pku.edu.cn

Our approach is generally a segmentation and classification framework. It takes a sequence of laser scan lines

as input that are represented as a range image, creates labeled segments as output, which can be converted to 3D point clouds of objects. In brief, the system contains three steps as shown in Fig.3. First, a series of pre-processing operations are conducted to filter out isolated points, sky and ground. Second, a graph-based segmentation is applied to find data clusters (called segments) corresponding to objects. Third, features are extracted to describe shape and spatial information of the segments, which are used in training a classifier.

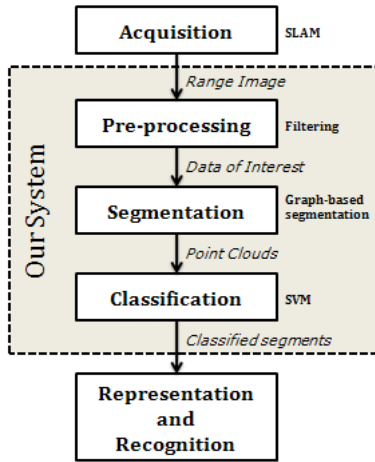


Fig. 3. Processing flow of the system

This paper is organized as follows. Section 2 gives a literature review on segmentation and classification methods of both point cloud and range image. Section 3 describes a graph-based segmentation. Section 4 addresses feature extraction and training classifier. Experimental results are presented in section 5, followed by conclusions and future studies in section 6.

II. LITERATURE REVIEW

Laser scan data, as a sequence of 2D or 3D points, represents environmental geometry directly. With its scanners becoming a standard equipment for a mobile robot, many researchers utilize its geometric advantage to assist scene understanding.

In practical, laser scan data is converted into a lower dimension, since it distributes sparsely in 3D space. In DARPA Urban Challenge, many teams use a $2\frac{1}{2}$ occupancy grid map to project 3D points to a horizontal plane [5], [14], where segmentation and classification of objects are done for higher-level reasoning. By fitting bounding boxes, the objects within a certain range, e.g. road boundaries can be detected, while these are restricted in a pre-defined environment. In [11], Himmelsbach et. al. demonstrate a mobile system using down sampling method and occupancy grids to classify objects in real-time.

In contrast, several papers label every 3D point directly. Based on examples, Angelov et. al. [1] and Lalonde et. al. [16] describe methods where every single point of a scan is assigned with a class label, considering a point and its neighbors are dependent. In this approach, labeling a point is influenced by labels in its local vicinity. Thus, Markov Random Fields are used to model their relationship. Douillard et. al. [3] propose a ground model and an object model to label semantic content in the urban scene. In [9], Golovinskiy et. al. investigate the design of a system that recognizes objects in 3D point clouds of urban environments based on shape descriptors and contextual information.

On the other hand, laser scan data can be represented in the form of a range image. As the format of a range image is consistent to a visual image, many methods developed in computer vision are of great reference. There is a large body of work addressing range image segmentation. A very famous report comparing the major segmentation methods can be found in [12]. Many of the methods are motivated by the needs for recognizing industry parts [15] or registering the data taken at different locations [25]. These works always assume simple or well-defined object geometry. There are still a few research works processing range images of real world scenes. [10] considers a real-world indoor and outdoor scene by modeling the man-made objects using planes and conics, free-form objects using splines, and trees using 3D histogram, segmentation and model fitting for each segment is formulated in a data-driven Markov Chain Monte Carlo procedure.

Motivated by the need of generating a semantic map of a large urban outdoor environment, where the environment is explored and sensed using a robot car with laser range scanners. We need to consider a scene that contains many kinds of objects, such as buildings, roads, trees, bushes, people, cars, etc., which have different scales in 3D space, with different geometric models. We refer to the researches in [17], [20], [24] that generate unified frameworks for the segmentation and recognition problems in a complex scene, comprising a mixture of objects.

III. SEGMENTATION OF A RANGE IMAGE

A. Pre-processing

As a single-row laser scanner measures the environment in a mode of scan-line by scan-line, range image is acquired where its horizontal axis indicates time, vertical axis indicates the sequential order of measurements and the depth is shown in grey scale.

First, isolated points are filtered out as noises. In a typical urban environment, we assume the ground is flat so that ground samples can be removed by plane fitting. Then points close to ground can also be grouped into ground in addition.

At the mean time, sky and glasses are also filtered out before segmentation is conducted. Since laser beam does not reflect heading for sky, there is a reading indicating infinity.

B. Graph-based Segmentation

Once the sky and ground are separated from the range image, the rest need to be segmented into objects we are interested in. As mounted on a mobile vehicle, data acquired by laser scanner usually spatially connect within its neighbors in a range image. This property suggests that we can use graph-based grouping techniques to accomplish the segmentation.

Our approach is most related to the graph theoretic formulation of grouping. The set of pixels is represented as a weighted undirected graph $G = (V, E)$, where V is the set of nodes standing for pixels and E is the set of edges between pixels and their neighbors in 3D space. There are many popular algorithms that cut this kind of graph, such as Normalized Cuts (NCuts)[22], the Felzenszwalb and Huttenlocher (FH) algorithm [4]. In our approach, we choose the FH algorithm, because it catches the non-local properties well and works efficiently.

The algorithm is outlined in Algorithm 1. Details about the internal difference and proof about the greedy property of this algorithm are discussed elaborately in [4]. The method runs in $O(m \log m)$ time for m edges and is also fast in processing our data.

Algorithm 1 Outline of the FH algorithm.

Input:

The undirected graph $G = (V, E)$;

Output:

The segmentation $S = (C_1, \dots, C_r)$ with r components.

- 1: Sort the edge set by weight in ascent order into (e_1, \dots, e_m) .
 - 2: Start with a segmentation $S^0 = (C_1^0, \dots, C_n^0)$, where each vertex is in its own component.
 - 3: **for** each edge e_i **do**
 - 4: Construct S^i given S^{i-1} . Let v_i^a, v_i^b denote the vertices e_i connects. If v_i^a, v_i^b are in disjoint components of S^{i-1} and the weight of e_i is less than the internal difference of both components, merge the two components otherwise do nothing.
 - 5: **end for**
 - 6: **return** $S = S^m$
-

C. Implementation Details

We first find k nearest neighbors for each pixel in the range image. The weight of the edges is simply assigned by Euclidean distance, i.e., pixels are grouped by spatial connectivity. There are two parameters, π and k , in the algorithm. Generally, they are related to the scale of observation. The larger k is, the larger component is preferred. In our mobile platform, data is always acquired along the road, then k is dependent on the distance between objects and scanners. Thus, we use a function $k = f(C)$ to formulate such dependence,

$$k = f(C) = \theta \cdot \text{Dist}^2(P_{CoG}, P_{Laser}) \quad (1)$$

where C is the component in segmentation S , P_{CoG} is the center of gravity of C , P_{Laser} is the position of laser scanner and θ is a pre-defined parameter. This can be calculated iteratively when a new pixel is added, so the time cost can be ignored compared to the whole segmentation algorithm.

Because the neighbors in image is also neighbors in 3D space, segmentation is done in a scan-line based procedure in our approach. As a result, we obtain segments representing objects of interest for further classification.

IV. CLASSIFICATION

A. Feature Extraction

Through a graph-based segmentation, the segments are represented as a set of components $S = (C_1, \dots, C_r)$. These components are viewed as point clouds in 3D space. Then we extract features from these point clouds for future classification.

Let $f_i^{(d_i)}$ denotes the i th d_i -dimensional feature of a point cloud. The feature set $f = \{f_1^{(d_1)}, \dots, f_M^{(d_M)}\}$ is used for classification based on point clouds. The features defined in this research are listed in Table I, where we need to emphasize the following two points.

1) *Normal vector estimation:* We use a fixed number of Euclidean nearest neighboring data points, say k points, to estimate the normal vector at a given data point. Therefore, a local polygonal mesh description of one point with its neighbors is created for further estimation. There are many approaches to estimate, such as plane-fitting, quadric surface fitting, area-weighted average method, angle-weighted average method, etc [19]. In our approach, average methods do not work better for arbitrary 3D data than surface fitting methods, as the quality of 3D mesh acquired in real world is not so good as that of object models for computer graphics due to noise and sampling rate. Finally we choose plane fitting because of its computation efficiency and effectiveness in describing the variance of a tree segment.

2) *Statistical features for components:* In order to classify the point cloud as a whole, we use statistical method to describe the points of an object. Both spatial and shape descriptors are considered here. For each point p_i in C , the properties can be acquired as $\{x_i, y_i, z_i, \mathbf{n}_i\}$, where (x_i, y_i, z_i) is its 3D coordinate and \mathbf{n}_i is its normal vector, calculated from its neighborhood of at most 24 points (5×5 excluding itself).

Furthermore, inspired by the spin image descriptor [13], we convert the normal vector of a point into measurement of longitude and latitude (we call it LLMap, see Fig.4), as a normal vector has only two degrees of freedom. We introduce a 2D histogram over this distribution. The prominent peaks correspond to the prominent surface in different directions, and their heights correspond to the saliency of surface respectively. These features are rotation invariant and effective in describing the shape of a object.

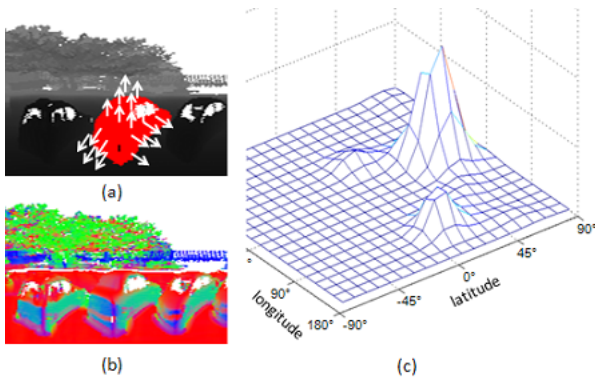


Fig. 4. (a) the original range image, white arrows show the direction of normal vector; (b) directions of the corresponding normal vector in RGB; (c) a histogram called LLMap, the two axes of horizontal plane are longitude and latitude respectively.

Feature	Definition
f_1	Horizontal size
f_2	Point deviation
f_3	Maximal height value
f_4	Z factor of Center of Gravity(CoG)
f_5	Variance of normals
f_6	Number of prominent peaks of LLMap
$f_7^{(6)}$	Statistics on LLMap

TABLE I
FEATURES EXTRACTED FROM A DATA SEGMENT

As a result, we combine common features and our features above to classify a point cloud. The feature set, as listed in Table.1, is selected into 12 dimensions after feature analysis. The likelihood (See Fig.5), measures $p(f_i^{(d_i)} | L_k)$ of a certain pair of feature $f_i^{(d_i)}$ and object label L_k .

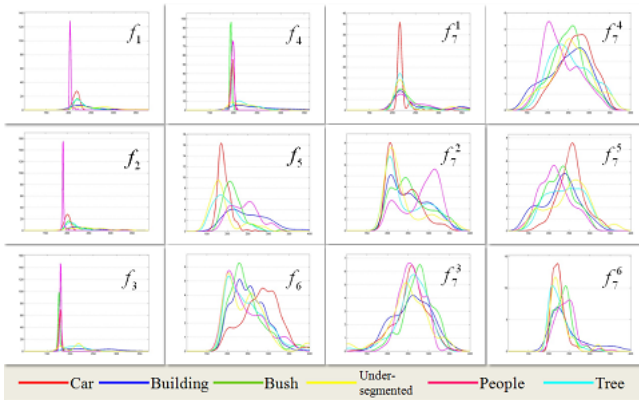


Fig. 5. Likelihood measures for classifying a point cloud. The definitions of these features are listed in Table.1.

B. Definition of the Classifiers

For classification of a segment, we consider both its local properties and evaluation of the whole segment as a point cloud. Let L denotes the set of object classes, i.e., $L =$

$\{building, bush, car, tree, pedestrian, bicycle, \dots\}$ and L_k indicates k th class label in the set. Thus, we formulate the classifier $p(C = L_k)$, which estimates the probability of the segment C belonging to label L_k .

1) *The SVM Classifier:* For the estimation of $p(C = L_k)$, a SVM classifier is chosen because our training set is much smaller than that of training line segments in our previous work. In this research, we refer to the off-the-shelf LibSVM[2] to give a probabilistic prediction. Cross validation is used for parameter optimization. We finally selected RBF kernel and set $C = 5, \sigma = 0.5$ by grid search.

2) *The Decision Tree Classifier:* To investigate the sensitivity of features, we also train a Decision Tree classifier. In this research, Classification and Regression Trees (CART) is used.

V. EXPERIMENTAL RESULTS

A. Training and Testing Samples

Data are collected inside the campus of Peking University using the system. The vehicle ran around a large building for a number of times at different situations. Fig.6 shows an example of integrated laser points used in this research. The data acquired by L5 (dark red) is manually segmented and labeled, while the data acquired by L4 (dark blue) is used for testing.

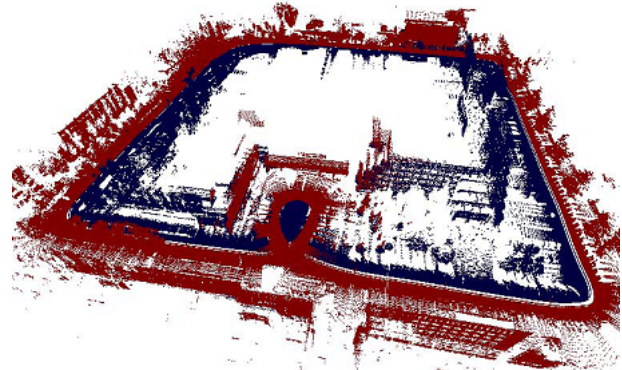


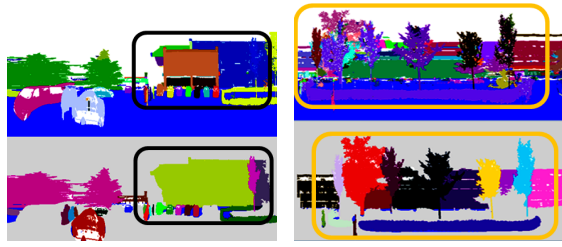
Fig. 6. A 3D view of laser data acquired by the laser scanners L4 and L5 of POSS-v.

Then segmentation and classification is evaluated separately.

B. Evaluation of Segmentation

We develop an interactive tool to make Ground Truth (GT) set manually. Based-on Random Walks [7], we can easily get segments by setting seeds. As in [12], we evaluate the automatic segmentation results by comparing with segments in GT set by pixels. In general, we consider 3 circumstances: under-segmented, over-segmented and correct-detected. We show two extreme cases in our research in Fig.7.

In under-segmented case, the classification could not continue since two or more objects are contained in the



(a) Over-segmented in Black Rectangle (b) Under-segmented in Orange Rectangle

Fig. 7. Extreme cases of the FH algorithm results

same segment. Thus, spatial connectivity is not enough for segmentation and more information of the object should be considered.

C. Evaluation of Classification

To evaluate the performance of two classifiers, we use traditional training and testing method. Both training and testing samples (See Table.2) are generated from automatic segmentation results after small segments are filtered out. From right part of Table.2, we may find both classifiers do well in classification of cars. The poor results for building suggest that multi-scale is an issue because a building is always over-segmented into pieces, and thus we may fail to obtain its overall properties.

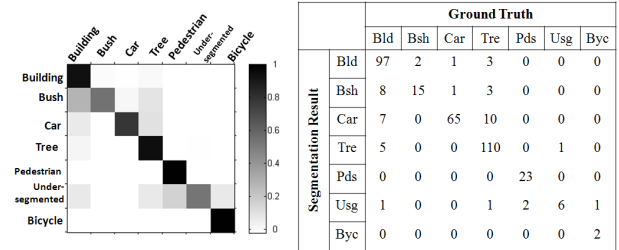
Object of Interest	Training Samples	Testing Samples	Precision of LibSVM	Precision of CART
Building	264	118	82.20%	87.45%
Bush	69	17	88.24%	94.12%
Car	216	67	97.01%	97.01%
Tree	364	127	88.19%	87.96%
Pedestrian	104	25	80.00%	88.00%
Bicycle	9	3	33.33%	100.00%
Under-segmented	24	7	85.71%	85.71%
Total	1050	364	86.81%	89.56%

TABLE II
TRAINING AND TESTING SAMPLES

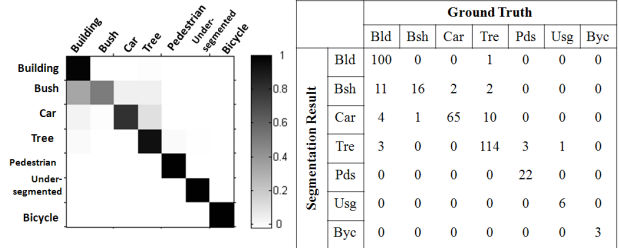
Fig.8 shows the confusion matrix for the data. The under-segmented class is special for the failure cases when two objects are too close to each other. Under these circumstances, we suggest that a recognition should be introduced to separate them. At the mean time, the normals of points in the segments provide 7 dimensions of features, only 3 of which are used in decision tree after pruning in our research. The superior performance of decision tree indicates that rule-based system rather than more features may work well as long as discriminative features are proposed.

D. Result

Using the method developed in this research, range image is partitioned into segments, meanwhile, labels representing



(a) SVM



(b) CART

Fig. 8. The confusion matrix of two classifiers

object types are associated to each individual segment. We present a clip of our results shown in Fig.9, and whole testing results in 3D view shown in Fig.10.

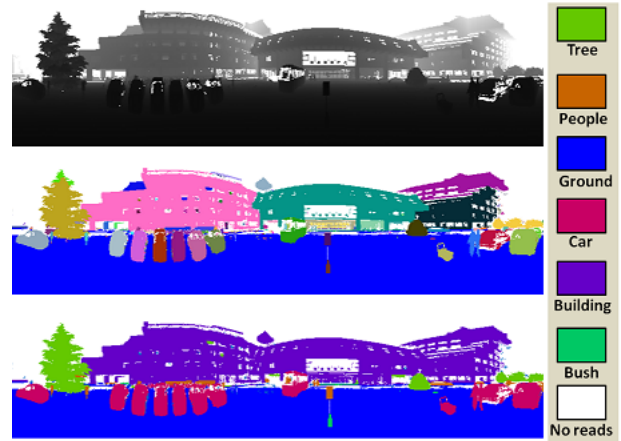


Fig. 9. A segmentation and classification result. Top: range image; Middle: segmentation result; Bottom: classification result.

VI. CONCLUSION AND FUTURE STUDIES

A. Conclusion

Given a sequence of laser scan measurements to the environment, this paper propose a method of segmentation and classification on range image, where 3D coordinates are also retrieved in calculation. The objective is to find data segments corresponding to objects and annotate class labels to them. A graph-based segmentation is applied to separate the data to different objects, and a classification is designed by extracting both local and global statistical features of

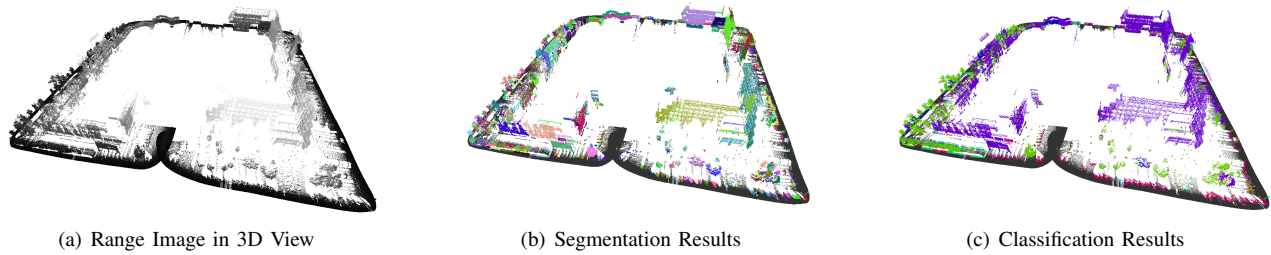


Fig. 10. 3D View of Testing Results

each data segment, and trained on manually extracted sample data. Efficiency of the method is demonstrated through experiments on the data of a complex urban outdoor scene. All the data used in this research, including both raw data and manually extracted training samples, are open freely at our website <http://poss.pku.edu.cn>.

B. Future Studies

In order to understand a complex scene, extending the object class is required, so as to contain the small scale objects, such as traffic signs, signal lights, guardrails, fences, etc. However, a key issue has to be addressed on generating training samples, which is quite label intensive, especially towards a study of real world scene. A method is to be developed of robustness to partial observations, while without relying on a large amount of training samples. In addition, as our final goal is to study scene semantic and support robot decision-making, a comprehensive investigate at different scenery is to be addressed.

REFERENCES

- [1] D. Anguelov, B. Taskar, V. Chatalbashev, D. Koller, D. Gupta, G. Heitz, and A. Ng, "Discriminative Learning of Markov Random Fields for segmentation of 3D scan data," *IEEE Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2005, pp 169-176.
- [2] C. Chang and C. Lin, LIBSVM: a library for support vector machines, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [3] B. Douillard, A. Brooks and F.T. Ramos, "A 3D Laser and Vision Based Classifier," *Proc. 5th Int. Conf. on Intelligent Sensors, Sensor Networks and Information Processing (ISSNIP)*, 2009
- [4] P.F. Felzenszwalb and D.P. Huttenlocher, "Efficient graph-based image segmentation," *Int. J. Computer Vision*, vol. 59, 2004, pp 167-181.
- [5] D. Ferguson, M. Darms, C. Urmson, and S. Kolski, "Detection, Prediction, and Avoidance of Dynamic Obstacles in Urban Environments," *IEEE Intelligent Vehicles Symposium*, vols. 1-3, 2008, pp 534-539.
- [6] C. Fruh and A. Zakhor, "3D model generation for cities using aerial photographs and ground level laser scans," *IEEE Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2001, pp 31-38.
- [7] L. Grady, "Random Walks for Image Segmentation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 11, 2006, pp 1-17.
- [8] A. Georgiev and P.K. Allen, "Localization methods for a mobile robot in urban environments," *IEEE Trans. on Robotics and Automation*, vol. 20, Oct. 2004, pp 851-864.
- [9] A. Golovinskiy, V.G. Kim, and T. Funkhouser, "Shape-based Recognition of 3D Point Clouds in Urban Environments," *IEEE Int. Conf. on Computer Vision*, 2009.
- [10] F. Han, Z. Tu, and S. Zhu, "Range image segmentation by an effective jump-diffusion method," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, 2004, pp 1138-1153.
- [11] M. Himmelsbach, T. Luettel, and H. Wuensche, "Real-time object classification in 3D point clouds using point feature histograms," *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2009, pp 994-1000.
- [12] A. Hoover, G. Jean-Baptiste, X. Jiang, P. Flynn, H. Bunke, D. Goldgof, K. Bowyer, D. Eggert, A. Fitzgibbon, and R. Fisher, "An experimental comparison of range image segmentation algorithms," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 18, 1996, pp 673-689.
- [13] A.E. Johnson, "Spin-images: A Representation for 3-d surface matching," *PhD Thesis*, 1997
- [14] S. Kammel, J. Ziegler, B. Pitzer, M. Werling, T. Gindele, D. Jagzent, J. Schöder, M. Thuy, M. Goebel, F.V. Hundelshausen, O. Pink, C. Frese, and C. Stiller, "Team AnnieWAYS Autonomous System for the DARPA Urban Challenge 2007," *The DARPA Urban Challenge*, 2009, pp 359-391.
- [15] D. Katsoulas, C.C. Bastidas, and D. Kosmopoulos, "Superquadric Segmentation in Range Images via Fusion of Region and Boundary Information," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, 2008, pp. 781-795.
- [16] J.F. Lalonde, N. Vandapel, D.F. Huber, and M. Hebert, "Natural terrain classification using three-dimensional lidar data for ground robot mobility," *J. Field Robotics*, vol. 23, 2006, pp 839-861.
- [17] T. Malisiewicz and A. Efros, "Recognition by association via learning per-exemplar distances," *IEEE Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2008, pp 1-8.
- [18] A. Nchter, H. Surmann, and J. Hertzberg, "Planning Robot Motion for 3D Digitalization of Indoor Environments," *Proc. 11th Int. Conf. on Advanced Robotics (ICAR)*, 2003, pp 222-227.
- [19] D. OuYang and H.Y. Feng, "On the normal vector estimation for point cloud data from smooth surfaces," *J. Computer-Aided Design*, vol. 37, 2005, pp 1071-1079.
- [20] J. Porway, K. Wang, B. Yao, and S.C. Zhu, "A hierarchical and contextual model for aerial image understanding," *Proc. Computer Vision and Pattern Recognition (CVPR)*, 2008, pp 1-8.
- [21] I. Posner, D. Schroeter, and P. Newman, "Online generation of scene descriptions in urban environments," *J. Robotics and Autonomous Systems*, vol. 56, Nov. 2008, pp 901-914.
- [22] J. Shi and J. Malik, "Normalized Cuts and Image Segmentation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, 2000, pp 888-905.
- [23] S. Thrun, W. Burgard, and D. Fox, "A real-time algorithm for mobile robot mapping with applications to multi-robot and 3D mapping," *IEEE Int. Conf. on Robotics and Automation (ICRA)*, vol.1, 2000, pp 321-328.
- [24] Z. Tu, X. Chen, A. Yuille, and S. Zhu, "Image parsing: unifying segmentation, detection, and recognition," *IEEE Int. Conf. on Computer Vision*, vol.1, 2003, pp 18-25.
- [25] J. Weingarten and R. Siegwart, "EKF-based 3D SLAM for structured environment reconstruction," *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2005, pp 3834-3839.
- [26] H.J. Zhao, L. Xiong, Z.G. Jiao, J.S. Cui, H.B. Zha, and R. Shibasaki, "Sensor Alignment Towards an Omni-Directional Measurement using an Intelligent Vehicle," *2009 IEEE Intelligent Vehicles Symposium*, vols. 1-2, 2009, pp 292-298.
- [27] H. Zhao, M. Chiba, R. Shibasaki, X. Shao, J. Cui, and H. Zha, "SLAM in a dynamic large outdoor environment using a laser scanner," *IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2008, pp 1455-1462.