

Contents

1	Introduction	2
2	Mechanical design	2
3	Electronics	2
3.1	NUC	2
3.2	Arduino	2
3.3	IMU	2
3.4	IR sensors	2
3.4.1	Calibration	2
3.4.2	Filtering	2
4	Motion control	2
5	Software	3
5.1	Global picture	3
5.2	Odometry	3
5.3	Mapping	3
5.4	Exploration	4
5.4.1	Wall following	4
5.4.2	BFS Search	4
5.5	Path planning	4
5.5.1	Global	4
5.5.2	Local	4
5.6	Localization	4
6	Computer Vision	4
6.1	Camera extrinsic calibration	4
6.1.1	Tilt compensation	5
6.2	Object detection	5
6.3	Object recognition	6
6.3.1	3D object recognition	6
6.3.2	2D object recognition	7
6.3.3	Final approach	7
6.4	Obstacle detection	8
7	Discussion	9
7.1	Results	9
7.2	System Evaluation	9
7.3	Project management	9
8	Conclusions	9

1 Introduction

This is the project report in course Robotics and Autonomous systems. The purpose of this report is to reflect ideas and solutions we provided for trying to reach the goal of this course – making autonomous robot that can explore maze, create map of it and detect and recognize different geometrical objects located in maze.

2 Mechanical design

Robot size is 23 cm x 23 cm width and length (including wheels) and 29 cm height. Same width and length were chosen due to ease of maneuverability. For the same reason wheels are located in the middle of its sides and corners of the robot are cut off in 45° angle. Two aluminium plates are connected with pillars in such way making two floors. On the first floor are located motors, accumulator battery and speaker whereas on the second floor is located NUC and Arduino sandwich. For camera there is pillar which raises it high enough so it reaches 29 cm height. IMU is located under the second plate as well as long range IR sensors which points to front and back direction. Short range IR sensors are located on the pillars which connect first and second floor.

$$y = f(x) \tag{1}$$

Equation 1 shows that...
Thrun *et al.*[1] show that...

3 Electronics

3.1 NUC

3.2 Arduino

3.3 IMU

3.4 IR sensors

3.4.1 Calibration

Because of non-linear sensor output vs distance curve it is not practical to use raw values from IR sensors in higher level nodes. Therefore it was needed to create node which converts raw data from sensors to distance. As the raw output changes values more rapidly in lower distances, measurements in low range were made quite densely (by measuring sensor output every 0.5 cm) whereas in longer distances such accuracy were not needed because output values were approximately the same even within 5 cm range. Sensor converter node is the only one which subscribes to raw sensors node and other nodes in ROS subscribes directly to distance data.

3.4.2 Filtering

To filter out noise from IR sensors data Kalman filter is used.

4 Motion control

To be able to reliably use motors instead of sending just PWM data to motors there needs to be a node which translates linear and angular speeds to necessary PWM values. However as PWM vs speed

relation is highly non-linear and depends on many factors such as ground friction it is not possible to just find suitable function for translating speed to PWM. Use of motor controller is therefore crucial. Its task is to control PWM values at the same time getting feedback from encoders. This data then can be used to adjust PWM values so that motor reaches the necessary speed and keeps it. In this project PID controller is used for motor control. Equation 2 shows general PID controller algorithm. In case of motor controller error $e(t)$ is difference between desired angular speed and current angular speed which can be calculated with equation 3. N describes ticks per revolution.

$$u(t) = K_p e(t) + K_i \int_0^t e(\tau) d\tau + K_d \frac{de(t)}{dt} \quad (2)$$

$$\omega(t) = \frac{2\pi \Delta encoder}{\Delta t * N} \quad (3)$$

With PID controller it is very important to find right parameters otherwise motion can be with oscillations, too slowly rising or could have other problems.

5 Software

In this project ROS software is used. It creates good and feature rich environment for simple nodes which can easily communicate to each other.

5.1 Global picture

5.2 Odometry

Odometry node takes raw data from motor encoders and translates it into global coordinates and angle according to equations 5.2. This data can be used for mapping and localization. Problem with odometry is that it tends to drift especially when speed increases.

$$distance_{left} = \frac{2\pi * Wheelradius * encoder_{left}}{ticksperrevolution} \quad (4)$$

$$distance_{right} = \frac{2\pi * Wheelradius * encoder_{right}}{ticksperrevolution} \quad (5)$$

$$\Delta x = \frac{distance_{left} + distance_{right}}{2} * \cos(\Phi) \quad (6)$$

$$\Delta y = \frac{distance_{left} + distance_{right}}{2} * \sin(\Phi) \quad (7)$$

$$\Delta \Phi = \frac{distance_{left} - distance_{right}}{Wheelbase} \quad (8)$$

5.3 Mapping

For map we chose to use occupancy grid package from ROS. IR sensors provide data about distances to the walls. This data together with position information from odometry node are used to make a map. Initially all area is unknown and by driving around and getting corresponding data from IR sensors free area and obstacles are detected and drawn in map. Data from sensors are quite noisy therefore some processing is needed to mitigate noisy results.

5.4 Exploration

For exploration there was created navigation node.

5.4.1 Wall following

5.4.2 BFS Search

5.5 Path planning

5.5.1 Global

5.5.2 Local

5.6 Localization

6 Computer Vision

Another big part of the project involves Computer Vision, allowing the robot to detect and identify objects in the maze. The task is subdivided into smaller subtasks that are now discussed.

6.1 Camera extrinsic calibration

First, a camera extrinsic calibration is required in order to be able to express measurements in camera frame (**camera-link**) into the **world** frame. In this process, a 6 DoF transform is computed. Given that the transformation between **world** and **robot** is computed by the **odometry** node, it is natural to compute here the transformation between **robot** and **camera-link**. Then, it is easy to just use the **tf** package to transform between coordinate frames.

In this particular case, we only compute 4 out of the 6 DoF of the pose, since we assume the roll and yaw of the camera to be 0. Therefore, we only compute the pitch (ϕ) and the translation vector \mathbf{t} .

Pitch

To estimate the pitch, we first take the Point Cloud published by the PrimeSense camera and extract the main plane, corresponding to the floor. Then, the normal vector \mathbf{n} is computed. Finally, we extract the pitch of the camera from Equation 9.

$$\phi = \cos^{-1} |n_y| \quad (9)$$

Translation To compute the translation vector $\mathbf{t} = \{t_x, t_y, t_z\}^T$, we take into account the transformation equation 10 in homogeneous coordinates:

$$\tilde{\mathbf{p}}^R = T_R^C \tilde{\mathbf{p}}^C \implies \begin{pmatrix} x^R \\ y^R \\ z^R \\ 1 \end{pmatrix} = \begin{pmatrix} R & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} x^C \\ y^C \\ z^C \\ 1 \end{pmatrix} \quad (10)$$

, where \mathbf{p}_R and \mathbf{p}_C are a 3D point in robot and camera coordinate frames, respectively, T_R^C is the transformation between the robot and the camera frame, R is the 3D rotation and \mathbf{t} is the translation vector.

R is computed from the yaw(θ), pitch (ϕ) and roll (ψ) angles:

$$R = R(\theta)R(\phi)R(\psi) = \begin{pmatrix} \cos(\phi) & 0 & \sin(\phi) \\ 0 & 1 & 0 \\ -\sin(\phi) & 0 & \cos(\phi) \end{pmatrix} \quad (11)$$

, since $\theta = \psi = 0$.

To compute the translation vector, we require only one known 3D point in robot coordinates that we try to transform into the camera coordinate frame. To do that, we use a "calibration sheet" as shown in Figure XXXXX.

We detect the four corners of the rectangle using a Fast Feature Detector from OpenCV and take the mass center as the known point. It is easy to extract the 3D coordinate by just getting it from the depth image. Finally, we compute the translation vector applying Equation 12.

$$\mathbf{t} = \mathbf{p}^R - R\mathbf{p}^C \quad (12)$$

This concludes the calibration, which provides a transformation from **robot** frame to **camera-link** frame, and is published to the **tf** tree afterwards.

6.1.1 Tilt compensation

Sometimes the robot tilts forward when braking and this greatly affect the calibration. Therefore, we implemented a small solution to fix this using the IMU. We basically compute the tilting of the robot (ϕ^R) taking the measurements from the accelerometer, a_x and a_z :

$$\phi^R = \tan^{-1} \left(\frac{|a_x - a_{x0}|}{a_z} \right) \quad (13)$$

Since the construction is not perfect, we first calibrate the IMU by storing the initial tilt, measured by a_{x0} .

Finally, the transformation between **robot** and **camera-link**, T_0 , computed before, is corrected according to Equation 14.

$$T' = T_{\text{IMU}} \cdot T_0 = \begin{pmatrix} R(\phi^R) & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} T_0 \quad (14)$$

6.2 Object detection

The core of the computer vision component is divided into two modules. Object detection is mean to be fast (close to real-time computing time) and robust. Therefore, our aim was to optimize it as much as possible. The object detection pipeline is presented in Figure XXXXXX.

1. First, we subscribe to the RGB and Depth images from the camera (for which we use the **ApproximateTime** and **Synchronizer** packages in ROS) and build a point cloud out of them. The reason for doing this is that for some reason the PrimeSense did not publish its own point cloud at a fixed rate, and it was quite far from 30 Hz (e.g.: every 200 ms or so), which was unacceptable. To make it faster, we downsampled the images by a factor of 4 in X and Y. For transforming points from 2D to 3D, we reverted the projection equation, as shown in Equation 15.

$$X = z \cdot \frac{u - c_x}{f_x} \quad ; \quad Y = z \cdot \frac{v - c_y}{f_y} \quad ; \quad Z = z \quad (15)$$

where $\{X, Y, Z\}$ is the 3D point in **camera-link** frame, $\{u, v\}$ is the 2D point from the RGB image, z is the depth at the given point, and c_x, c_y, f_x, f_y are the camera intrinsics, which are obtained from the **CameraInfo** message.

2. We transform the point cloud into the **robot** coordinate frame and extract the floor as a separate point cloud. This is done using a simple **PassThrough** filter, preserving the points whose **z** component is smaller than 0.01 m. We could have used the **Plane Segmentation** library from PCL, but this was more computationally expensive and would not always work since it would remove only the dominant plane, which might not be always the floor.
3. A binary mask is created representing the floor by just reprojecting the floor point cloud extracted before. The process is done simply by reverting Equation 15 and getting u and v . Since a subsampling was done at the beginning, it was necessary to apply dilation in order to have a solid mask.
4. Next, the floor is removed from the original RGB image by an AND operation with the negated floor mask. Now it that the yellow floor will not trigger false positives, it is possible to perform **color filtering**. We select to do this in the HSV color space given its better robustness against illumination. This is performed using a bank of 5 color filters (red, green, blue, yellow, purple) with ranges in H and S manually tuned for the application. The function **inRange** of OpenCV is used to quickly perform this filtering. The result is a set of binary masks, from which we select the one with a larger color response. After that, we find contours and filter them by size (a minimum is required) and the aspect ratio (a maximum value of 2 is used).
If a contour still remains, it is likely that it belongs to an object. The biggest contour that fulfill these conditions is taken and a binary mask is created out of it using the **drawContours** function.
5. In case a contour belonging to a coloured object was found, we compute the 3D position of its mass center both in robot and world coordinates. If the object has not been recognized yet, and if it is close enough to the robot (less than 30 cm from it), the recognition module is called.

XXXXXX Images from the process

6.3 Object recognition

A different number of approaches for object recognition have been tried. We decided to have a parallel research in both 2D and 3D methods for object recognition to analyze their advantages and disadvantages.

6.3.1 3D object recognition

First, several attempts to perform 3D recognition were tried. We started using **feature matching** object recognition. The first question to ask was: what features? The task objects were too simple and had not texture, and therefore no keypoints could be extracted from them. Therefore, we tried to just subsample the point cloud to see if those points could act as keypoints.

Then, we tried several feature descriptors. We started with SHOT and PFHRGB, which were proven to provide a great accuracy (XXXXXX REF). They did work well for cubes and balls, but we sadly discovered that it was not possible to do that with the hollow objects: they became "invisible" to the PrimeSense due to interferences with the IR projection, as can be seen on Figure XXXXXX:

The computational time was also quite expensive (around 1-2 seconds), and it was not robust against illumination since it was using RGB data. We further tried shape-only descriptors (PFH and FPFH) without a significant increase in performance. Therefore we abandoned the idea of 3D-only recognition.

6.3.2 2D object recognition

RYAN

6.3.3 Final approach

The final approach was to determine a **probabilistic formulation** of the problem in order to gain in robustness and be able to reason about the certainty of the recognition. The most likely object o^* was determined using a MAP estimate given the measurements in 3D shape z^s and color z^c , according to Equation 16.

$$o^* = \arg \max_o p(o|z^s, z^c) = \arg \max_o \eta p(z^s, z^c|o) p(o) = \arg \max_o p(z^s|o) p(z^c|o) \quad (16)$$

Here we are assuming that all the objects have the same prior probabilities and that the shape and color are conditionally independent given the object. We now describe the way to compute the previous probabilities.

3D shape probability For the task of 3D shape recognition, we make use of the VFH descriptor included in the PCL library. Unlike the previously described descriptors (FPFH, PFHRGB,...), this is a **global** descriptor, which means it does not have to be applied to all possible keypoints of the object, but rather it is applied to the whole point cloud. It is rotation and scale invariant, which is important for the application.

We build a 3D model database for the objects. We realized early in the course that **hollow objects** are invisible to the PrimeSense due to interference generated by the IR projector. Therefore, we only built 3D models for the ball and cube. 50 models were taken for each of them at different distances from the camera and with different poses.

After that, the test point cloud (which contains a potential object) is matched against the database. This is done quickly by histogram matching, for which OpenCV has some metrics. We opted for using the **correlation** metric between two histograms H_1 and H_2 , according to Equation ??

$$d(H_1, H_2) = \frac{\sum_i (H_1(i) - \bar{H}_1(i))(H_2(i) - \bar{H}_2(i))}{\sqrt{\sum_i (H_1(i) - \bar{H}_1(i))^2 \sum_i (H_2(i) - \bar{H}_2(i))^2}} \quad ; \quad \bar{H}_k = \frac{1}{N} \sum_i H_k(i) \quad (17)$$

Figure XXXX shows the histograms returned by the VFH descriptor. The difference between a cube and a ball is quite clear.

A value of -1 means complete opposite correlation, while 1 means perfect match. Therefore, we compute the probability of 3D shape as follows:

$$p(z^s = \text{cube}|o) = \max_i 0.5 \cdot (d(H_{\text{test}}, H_{\text{cube}}^i) + 1.0) \quad (18)$$

$$p(z^s = \text{ball}|o) = \max_i 0.5 \cdot (d(H_{\text{test}}, H_{\text{ball}}^i) + 1.0) \quad (19)$$

$$(20)$$

where H_{cube}^i and H_{ball}^i are the descriptors of the 3D models from the database.

Hollow probability

We also take into account the probability that an object is hollow, according to Equation 21.

$$p(z^s = \text{other}|o) = 1.0 - N_d/N_T \quad (21)$$

, where N_d is the number of pixels inside the object mask that have a valid (non-zero, not NaN) depth and N_c is the number of pixels within the object mask.

These three probabilities are normalized afterwards so that they add up to 1.

Color probability

In this task, it is not possible to recognize all the objects without color information. For this we implemented a Bayes Classifier XXXX ref on the HSV color model. We currently only use the H and S components, which influence most the outcome.

The first step is to create a **color model**. We decided to have independent Gaussian distributions for 7 colors (red, green, blue, yellow, purple, orange, light green) and for each of the channels (H and S), so the total number of parameters is $7 \cdot 2 \cdot 2 = 28$. For every color, the model $\mu_h, \sigma_h, \mu_s, \sigma_s$ is computed from training data. In particular, 50 pictures from each of the 10 objects were taken at different positions and different light conditions. In sum, 500 pictures form the color model. An application was implemented in order to speed up the capture of data. In particular, the application would cut out an image from a predefined region of interested at a fixed rate (e.g.: every 5 seconds). Finally, the color model was computed by extracting the mean and variance values for each color and channel.

During the test phase, the Bayes Classifier was applied to every pixel included in the object mask to determine to which class it belonged to. Finally, the probability of any color is computed as follows:

$$p(z^c = c|o) = N_c/N_T \quad (22)$$

, where N_c is the number of pixels belonging to class c , and N_T is the total number of pixels in the object mask.

6.4 Obstacle detection

In addition to the so-called "IR bumper" (see section XXXXXXXXXXXXXXXX), an obstacle avoidance solution using the PrimeSense camera was implemented. The aim was to detect also objects right in front of the robot, especially the edge of walls that would not be detected by the front IR sensor.

To do this, we simulated a laser scanner using the depth image. The implemented pipeline is as follows:

1. Build a point cloud from the depth image and transform it into the **robot** coordinate frame. Here color was not required.
2. Filter the point cloud to determine the working region. In particular, the 3D points following these constraints were selected:
 - $0.2 < x < 0.4$
 - $-0.115 < y < 0.115$
 - $0.05 < z < 0.06$

All the values are expressed in meters. In other words, a thin slice (0.01 m thick) covering an area of a width equals to the robot width (23 cm) and a maximum depth of 0.4 m (from the center of the robot).

3. Simulate laser beams. For this, the previous point cloud was filtered multiple times (once per line) using the constraint: $i < y < i + 0.005$, for $i = -0.115, \dots, 0.115$, increasing in steps of 0.005 meters. That is: lines with 0.005m resolution where taken.
4. Analyze the occupancy of the laser beam. This was done simply by checking if the previously filtered cloud lines contained any point. In that case, information about the minimum depth at which there was a point in the cloud was obtained.

The result of the algorithm is a set of lines in the forward direction of the robot that determine the occupancy of the region in front of it. An example is shown in Figure XXXXX

This information is afterwards sent to the `mapping` node, which sets to *blocked* the cells corresponding to the end point of the lines of the laser scan that do not show free space.

7 Discussion

7.1 Results

7.2 System Evaluation

7.3 Project management

8 Conclusions

Conclusions

References

- [1] S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005.