

UNIVERSITY OF SOUTHERN DENMARK

COURSE PROJECT

ROBOTICS 5<sup>TH</sup> SEMESTER - FALL 2018

---

## Marble finding robot

---



---

Alexander Tubæk Rasmussen  
alras16@student.sdu.dk

---

Kenni Nielsen  
kenil16@student.sdu.dk

---

Marcus Enghoff Hesselberg Grove  
grov16@student.sdu.dk

# 1 Abstract

# Contents

<b>1</b>	<b>Abstract</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>1</b>
2.1	Problem statement . . . . .	1
2.2	Readers guide . . . . .	1
<b>3</b>	<b>Design</b>	<b>2</b>
3.1	Environment . . . . .	2
3.2	Lidar Scanner . . . . .	2
3.3	Tangent bug . . . . .	6
3.4	Search strategy . . . . .	7
3.5	Q-learning . . . . .	7
<b>4</b>	<b>Implementation</b>	<b>9</b>
4.1	Line detection . . . . .	9
4.2	Tangent bug . . . . .	9
4.3	Search strategy . . . . .	9
4.4	Q-learning . . . . .	9
<b>5</b>	<b>Discussion</b>	<b>10</b>
<b>6</b>	<b>Conclusion</b>	<b>11</b>
	<b>Appendices</b>	<b>12</b>
<b>A</b>	<b>Tests</b>	<b>12</b>
A.1	Room based probability of marbles spawning . . . . .	12
A.2	Estimate for path lengths . . . . .	13
A.3	The impact of $\epsilon$ on Q-learning performance . . . . .	14
A.4	The impact of $\alpha$ on Q-learning performance . . . . .	16
A.5	Motion planners in action . . . . .	18

## 2 Introduction

### 2.1 Problem statement

#### 2.1.1 Problems

The following problems are stated to better describe the focus points throughout the project.

XXXX

1. How can an algorithm, which can find marbles in a closed map, be designed?
2. How can this algorithm be optimized, so the robot can find the ball faster, based on previous trails?
3. How can marbles and obstacles be detected from lidar data?
4. How can marbles and obstacles be detected from camera data?
5. How can a Fuzzy-control algorithm be constructed?
6. How can computer vision be used to detect marbles and obstacles, and to construct the map?

XXXX

#### 2.1.2 Limitations

The robot must be able to find marbles in a closed map and optimize itself based on previous trails. Thus the following limitations to the algorithms is listed:

1. To detect marbles and obstacles.
2. To construct a map.
3. To optimize for the shortest distance and time.

QT Creator, Fuzzy Lite, Gazebo and Matlab will be used throughout the project.

### 2.2 Readers guide

### 3 Design

#### 3.1 Environment

The two-wheeled robot should navigate around the environment called “bigworld” shown on figure 1a. To do so, it have been chosen to divide the environment into rooms. The natural definition of ‘rooms’ have been taken into account, resulting in a 14 room layout. This can be seen on figure 1b, where each room can be distinguished by different grey scale colours.

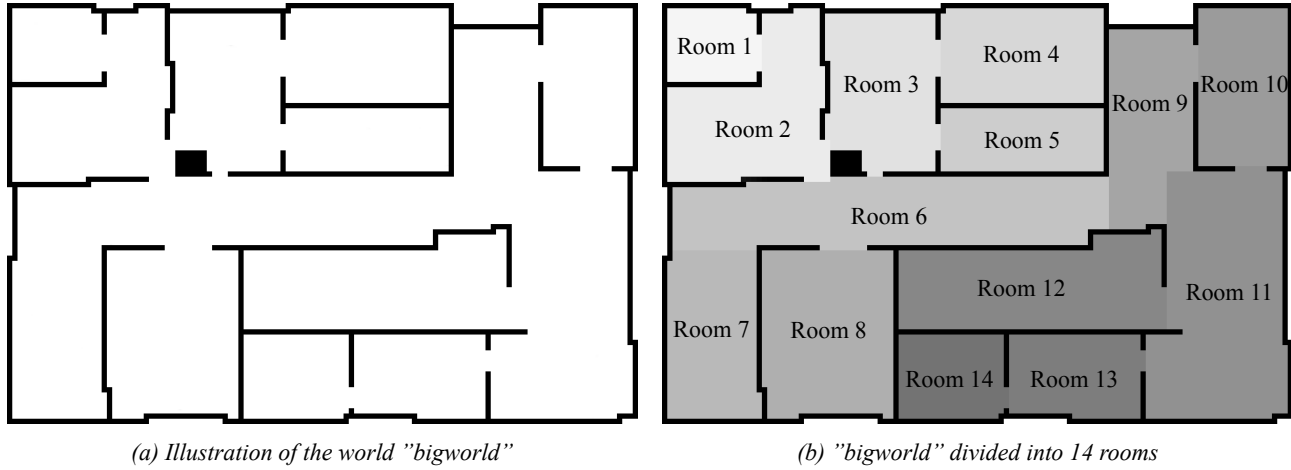


Figure 1: Illustration of the world “bigworld” before and after it has been divided into 14 rooms

The division into rooms are useful in terms of knowing which rooms that have been searched and which that not yet have been searched for marbles.

It is also useful as an abstract state space for reinforcement learning such that it can be found which order the rooms should be visited.

#### 3.2 Lidar Scanner

The two-wheeled robot given for this project is among other equipped with a 2d lidar (Light Detecting and Ranging) scanner. A lidar scanner detects the distance to targets by emitting a laser pulse and analyzing the time it takes for the beam to reflect and return to its source. The lidar scanner maps the environment and has a detecting range of 10 pixels in Gazebo and a 260 degrees field of view.

The script converts this range into mm by first scaling the pixel map to double size, such that the detecting range is 20 pixels. Afterwards, the script trace the pixel map to a eps-figure using 72 dpi as standard for the postscript. Then, the script converts this range from inches to mm by using a conversion factor of 25.4 mm per inch. This means that this range can be converted into mm by using the following formula:

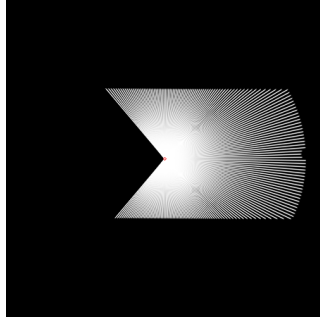
$$\frac{20 \text{ pixels}}{72 \text{ dpi}} \cdot 25.4 \frac{\text{mm}}{\text{inch}} = 7.06 \text{ mm}$$

Now the script scales the world from mm to m, which means that the range of the lidar scanner therefore is 7.06 m.

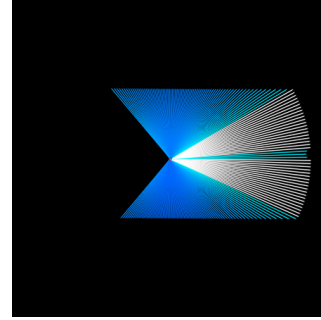
The lidar scanner maps the surrounding environment by collecting 200 datapoint, which must be processed in order to recognize objects such as walls and marbles. This means that circle and line detecting algorithms must be writing.

The datapoints from the lidar scanner is first visualized by drawing white lines from the robot’s location to each of the 200 datapoints using the `cv::line()` function as shown on figure 2a. Afterwards the datapoints are sorted by checking if the range is equal to max detecting range in order to get the points, that reflects from different object.

The filtered datapoints are drawn as lines upon the unsorted data. These lines can be distinguished by different blue colours depending on the range between the two points ranging from blue to cyan. This blue coloured and white lines are shown on figure 2b.



(a) Illustration of the unfiltered data



(b) Illustration of the filtered data

Figure 2: Illustration of the unfiltered and filtered data

The two sections below, describes the design of a line and a marble detection algorithm.

### 3.2.1 Line detection

It is usually important for a mobile robot to know its environment. There are several reasons for that, one is that robot must know the location of the obstacles (walls) relative to it to avoid driving into them. In the "bigworld" environment, the obstacles is walls which can be represented as lines from the filtered datapoints. These lines can be represented using a normal parametrization in polar coordinates, given by the following formula:

$$l : r = x \cdot \cos(\alpha) + y \cdot \sin(\alpha) \quad (3.1)$$

where  $r$  represents the distance from the origin to the closest point on the line and  $\alpha$  is the angle between the x-axis and the plane normal.

The Total Least Square method assumes that a line can be represented as in equation 3.1. This method also involves a determination of the orthogonal distance (the shortest distance) from a point  $p_i$  to a line  $l$  as shown on figure 3. The normal parametrization of the line  $l_i$  is given by the following formula:

$$l_i : r_i = x_i \cdot \cos(\alpha) + y_i \cdot \sin(\alpha) \quad (3.2)$$

The separation between those two lines ( $l$  and  $l_i$ ) is given by the difference  $d_i = r_i - r$ , since both lines have the same  $\alpha$ . This means that the orthogonal distance can be described using the following formula:

$$d_i = x_i \cdot \cos(\alpha) + y_i \cdot \sin(\alpha) - r \quad (3.3)$$

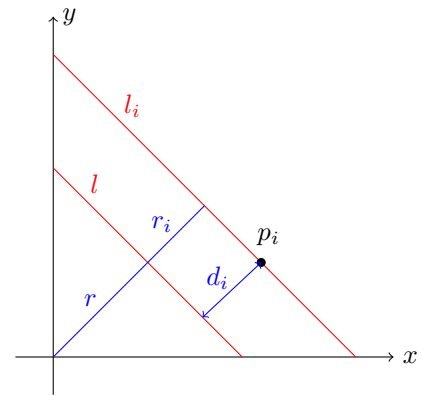


Figure 3: Orthogonal distance from point  $p_1$  to line  $l$

This only applies if we assume that there is no noise on the measurements.

This method gives an solution to the problem of fitting a straight line to a dataset of points  $p$  with  $n$  measurements having errors. The problem of fitting a line can be determined using the following sum:

$$\chi^2(l, z_1, \dots, z_n) = \sum_{i=1}^n \left[ \frac{(x_k - X_k)^2}{u_{x,k}^2} + \frac{(y_k - Y_k)^2}{u_{y,k}^2} \right] \quad (3.4)$$

where  $(x_k, y_k)$  are the points coordinates with corresponding uncertainties  $(u_{x,k}, u_{y,k})$  and  $(X_k, Y_k)$  denote its corresponding point of the straight line  $l$ . In the case of fitting the best line to the dataset, minimizes the expression for  $\chi^2$  by setting  $u_{x,k} = u_{y,k} = \sigma$  and  $k = 1, \dots, n$ . This reduces the problem to the Total Least Square method and minimizing is equal to minimizing the orthogonal distance of the measurements to the fitting line. Therefore, in the case of fitting the best line minimizes the expression above to the following:

$$\chi^2(l; Z) = \sum_{i=1}^n \frac{d_i^2}{\sigma^2} \quad (3.5)$$

$$= \sum_{i=1}^n \frac{(x_i \cos(\alpha) + y_i \sin(\alpha) - r)^2}{\sigma^2} \quad (3.6)$$

$$= \frac{1}{\sigma^2} \cdot \sum_{i=1}^n (x_i \cos(\alpha) + y_i \sin(\alpha) - r)^2 \quad (3.7)$$

A condition for minimizing  $\chi^2$  is done by solving the nonlinear equation system with respect to each of the two line parameters ( $r$  and  $\alpha$ )

$$\frac{\partial \chi^2}{\partial r} = 0 \quad \frac{\partial \chi^2}{\partial \alpha} = 0 \quad (3.8)$$

The solution of this nonlinear equation system is determined to the following:

$$r = \bar{X} \cos(\alpha) + \bar{Y} \sin(\alpha) \quad (3.9)$$

$$\alpha = \frac{1}{2} \arctan \left( \frac{-2 \sum_{i=1}^n [(x_i - \bar{X}) - (y_i - \bar{Y})]}{\sum_{i=1}^n [(x_i - \bar{X})^2 - (y_i - \bar{Y})^2]} \right) \quad (3.10)$$

where  $\bar{x}$  and  $\bar{y}$  are the means of  $x$  and  $y$ .

As explained earlier, the two-wheeled robot should avoid obstacles (walls). To do so, the robot needs to know the location of the walls. It is done by processing the datapoints from the lidar scanner and determining the points which fits to a straight line using a line extraction algorithm. There are several different line extraction algorithms to choose from. The incremental line extraction algorithm is chosen, because it is simple to implement. The pseudo code for the incremental algorithm is shown below.

#### Algorithm: Incremental

1. Start by the first 2 points, construct a line
2. Add the next point to the current line model
3. Recompute the line parameters
4. If it satisfies line condition (go to 2)
5. Otherwise, put back the last point, recompute the line parameters, return the line
6. Continue to the next 2 points, go to 2

This algorithm implements the Total Least Square method then computing the line parameters. Furthermore, the line model consists of points, which all must comply some line conditions. In that case all points must comply these conditions. The first condition is a threshold for the angle between the previous and current line model as shown on figure 4a. The threshold is defined to be the following:

$$\theta_{max} = 0,0025$$

The second condition is the angle between two points relative to the robot location as shown on figure 4b. This angle

should be greater than this difference, but not twice as great, since this condition should separate the points into two lines, if one point is missing on the list as shown on figure 4c. This angle is therefore defined to be the following:

$$\Delta\theta = (\theta_0 - \theta_1) \cdot 1,25$$

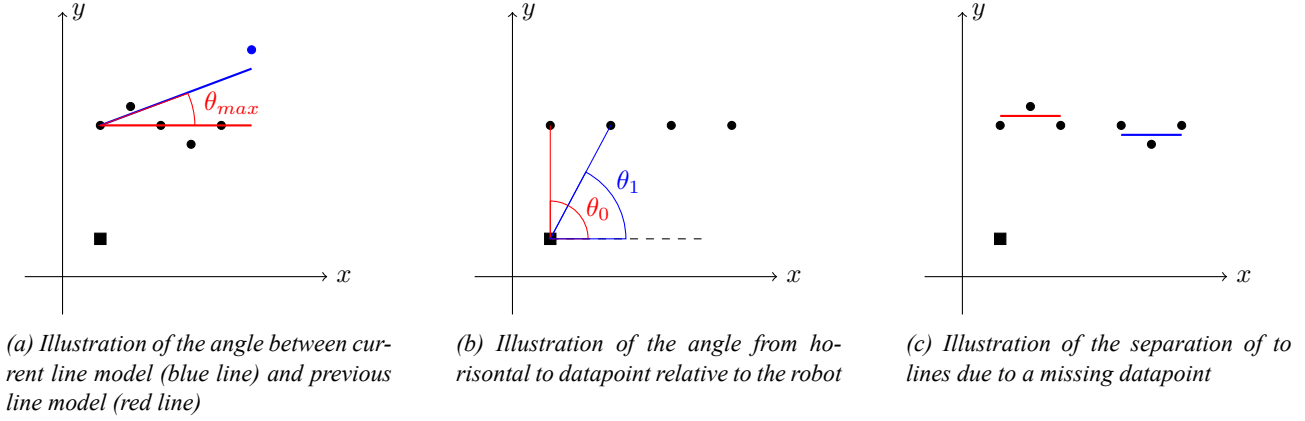


Figure 4: Illustration of the angle between two datapoints and two line models

### 3.2.2 Marble detection

The objective for the two-wheeled robot is to collect marbles effectively, while avoiding obstacles such as walls. The Hough transform is a algorithm, that maps from image space to the probability of the existence of features such as lines, circles or other general shapes. The algorithm is also capable of detecting partial object as in the case of the lidar scanner. However the algorithm are not very robust to noise, and can be very hard to calibrate in order to make the result from the algorithm useful. Furthermore the algorithm have a high computational cost.

Because of this, it have been chosen to design a simpler method for finding marbles in the lidar data.

Due to the fact that only a part of the circles would be visible from the data, it was chosen to calculate the center and radius of the circles from the chord and arch height of the circles. This can be used to determine the location of the marbles relative to the robot. Given this information the planner would be able to drive towards the marbles and "collect" them. It is assumed that the two detected outer points on the circle periphery defines the circle chord. It is also assumed that the orthogonal distance from the detected middle point on the circle periphery to the circle chord defines the circle arch height. This only applies for an uneven number of detected point.

The circle chord can be determined using the formula:

$$C = 2r \cdot \sin\left(\frac{\theta}{2}\right)$$

The circle camber can be determined using the formula:

$$h = r \left(1 - \cos\left(\frac{\theta}{2}\right)\right)$$

Solving this equation system, the circle radius was found to be:

$$r = \frac{K^2 + 4h^2}{8h}$$

The polar coordinates of the circle center can be found by adding the radius to the range of the detected middle point.



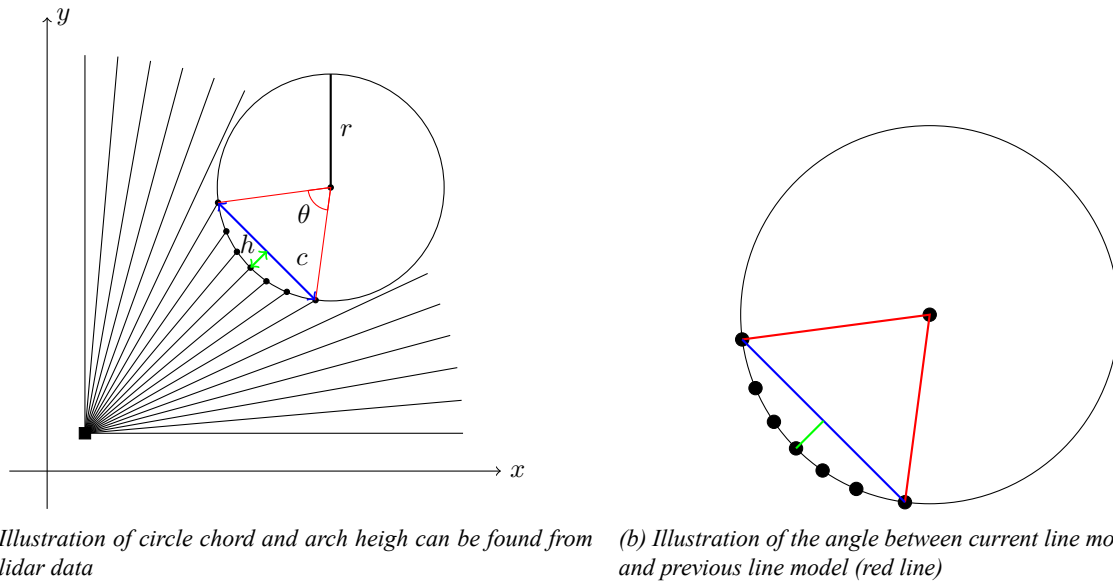


Figure 5: gegegegeger egeg gege

### 3.3 Tangent bug

The tangent bug algorithm is a sensor based planner which relies on inputs from a sensor to determine whether it should continue towards a specified location  $q_{goal}$  or to follow an obstacle  $O_i$  until a free path is available. The advantage of using this algorithm is the computational minimization and the fact that it only needs a start and end point. The reason for using this algorithm is the benefits just explained and that it will always try to achieve a shortest path solution to the goal location. Pseudo code for the implementation for the tangent bug algorithm is shown below.

#### Algorithm: Tangent Bug Algorithm

**Input:** A robot with a range sensor

**Output:** A path to the  $q_{goal}$  or a conclusion no such path exist

**while** True **do**

**repeat**

    Continuously move toward the point  $n \in (T, O_i)$  which minimizes  $d(x, n) + d(n, q_{goal})$

**until**

- The goal is encountered **or**
- The direction that minimizes  $d(x, n) + d(n, q_{goal})$  begins to increase  $d(x, q_{goal})$  so the robot detects a local minimum  $d(x, O_i) + d(O_i, q_{goal})$  on the boundary  $i$ .

Now choose the boundary following direction which continuous in the same direction as the most recent motion-to-goal direction.

**repeat**

    Continuously update  $d_{reach}$ ,  $d_{followed}$  and  $O_i$ .

    Continuously moves toward  $n \in \{O_i\}$  that is the chosen boundary direction.

**until**

- The goal is reached
- The robot completes a cycle around the obstacle in which case the goal cannot be reached.
- $d_{reach} < d_{followed}$

**end while**

It has been decided to make small changes to the algorithm and make it greedy. Instead of following the boundary  $d_{reach} < d_{followed}$ , the robot will follow that obstacle that minimizes  $d(x, n) + d(n, q_{goal})$ . The reason for doing this is to reduce the path from a giving point to a target location.

### 3.4 Search strategy

### 3.5 Q-learning

In order to effectively search the environment and collect marbles, a good search strategy must be found. This can be done by utilising reinforcement learning. By using reinforcement learning, the robot can learn from its experience and obtain a good strategy for navigating the environment.

By using a Temporal-Difference learning strategy the optimal action-value function can be estimated by every move taken unlike a Monte Carlo strategy where an episode terminates before any learning is obtained. In some cases with long episodes the Monte Carlo strategy is considered too slow.

Generally there are two categories of Temporal-Difference learning; on-policy and off-policy methods. One of the advantages of an off-policy over an on-policy method are that the action-value function can be estimated independent from the policy being used. The policy only influences which state-action pairs that are visited and updated.

Based on this Q-learning are chosen, the Q-learning method builds on the following update function for updating the action-value function (Q-values).

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)] \quad (3.11)$$

The update function for Q-learning consists of the old value for a given state-action combination plus a scaled difference between the old value, the immediate reward and the maximal value for the next state. The learning rate are denoted  $\alpha$  and ranging from 0-1 preferable closer to 0, in order to not to base the policy on this action only.  $\gamma$  denotes the discount factor and are also ranging from 0-1, preferable closer to 1, to ensure that future actions matter.

In the box below, the algorithm for Q-learning can be seen.

#### Algorithm: Q-learning

```

Algorithm parameters: step size:  $\alpha \in [0,1]$ , small  $\epsilon > 0$ 
Initialise  $Q(s, a)$ , for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
  Initialise  $S$ 
  Loop for each step of episode:
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon - greedy$ )
    Take action  $A$ , observe  $R, S'$ 
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal

```

#### 3.5.1 Definition of states

In order to perform Q-learning a definition of states must be made. These states must have the Markov property meaning that the probability of a marble being in a room must not depend on whether other rooms have been visited or not.

In order to achieve this, it have been chosen to implement a vector of boolean values, one element for each room. This will be used to store which rooms have been visited.

By doing this all possible combinations will be possible and the reward of entering a room will not depend on the other rooms, only the room itself.

It have also been chosen to implement the state with an integer storing the room number, and a boolean for storing whether the state is terminal or not. The definition of the state can be seen below.

```
qState
    int roomNumber
    std::vector<bool> roomsVisited
    bool isTerminal
```

## 4 Implementation

### 4.1 Line detection

#### 4.1.1 Marble detection

### 4.2 Tangent bug

### 4.3 Search strategy

### 4.4 Q-learning

Q-learning was implemented by designing a class containing the functionality of holding the state space and update the Q-matrix. This class was named `q_learning`.

Due to the definition of the state, each room will have  $2^{\text{number of rooms}}$  number of states, corresponding to the number of possible combinations of how the rooms could be visited. This will result in a very large Q-matrix. This could be overcome by splitting the Q-matrix into several smaller Q-matrices and then keep track of the mapping between them.

The Q-matrix would consist of number of base states by the number of actions which would be a 6 x 6 matrix with 5 rooms (including an initial state)

	Start	Room 1	Room 2	Room 3	Room 4	Room 5
Start	0	-100	-100	13.33971	-100	-100
Room 1	-100	0	8.681309	-100	-100	-100
Room 2	-100	0.405394	0	10.74771	-100	-100
Room 3	-100	-100	7.757309	0	7.837891	17.456
Room 4	-100	-100	-100	11.01171	0	-100
Room 5	-100	-100	-100	10.79571	-100	0

Table 1: Rewards for all state-action combinations given that no rooms have been visited and initial state is room 3

5 Discussion

## 6 Conclusion

# Appendices

## A Tests

### A.1 Room based probability of marbles spawning

The purpose of this test is to determine the probability of a marble spawning in each of the 14 rooms described in section 3.1.

#### A.1.1 Description of test

This test was done by conducting a total of 50 tests, where the position of the 20 marbles in the Gazebo environment are saved resulting in a total of 1000 samples.

These marbles were then mapped to one of the 14 rooms using the class `map_class`. The total amount of marbles found in each room can be seen in table 2. This data was then divided by the total number of marbles, to find the probability.

Due to the fact that the rooms are not the same size, this probability was divided by the size of the room (number of pixels in the map) and then normalised. The result can be seen in table 3.

#### A.1.2 Test parameters

- World used	bigworld
- Number of spawned marbles	20
- Number of tests	50

#### A.1.3 Data

Distribution of marbles	
Total number of marbles found in room 1	5
Total number of marbles found in room 2	65
Total number of marbles found in room 3	75
Total number of marbles found in room 4	52
Total number of marbles found in room 5	72
Total number of marbles found in room 6	158
Total number of marbles found in room 7	17
Total number of marbles found in room 8	118
Total number of marbles found in room 9	100
Total number of marbles found in room 10	22
Total number of marbles found in room 11	75
Total number of marbles found in room 12	160
Total number of marbles found in room 13	47
Total number of marbles found in room 14	34

Table 2: Table of how the marbles are distributed in the 14 rooms based on all 50 tests

Probability of marbles	
Probability of marbles found in room 1	0.012232
Probability of marbles found in room 2	0.061057
Probability of marbles found in room 3	0.078610
Probability of marbles found in room 4	0.059975
Probability of marbles found in room 5	0.117993
Probability of marbles found in room 6	0.098801
Probability of marbles found in room 7	0.019629
Probability of marbles found in room 8	0.099063
Probability of marbles found in room 9	0.114518
Probability of marbles found in room 10	0.027633
Probability of marbles found in room 11	0.044562
Probability of marbles found in room 12	0.130379
Probability of marbles found in room 13	0.072915
Probability of marbles found in room 14	0.062541

Table 3: Table of how the probabilities are distributed in the 14 rooms with room size taken into account

#### A.1.4 Conclusion

It can be concluded that the highest number of marbles was found in room 12 closely followed by 6 and 8. But due to the size of the rooms, the highest probability is found in room 12, followed by 5 and 9.

## A.2 Estimate for path lengths

The purpose of this test was to find an estimate for the distances between the rooms in order to have a distance punishment for the Q-learning.

### A.2.1 Description of test

This test was conducted using Geogebra a cas tool for geometry and algebra. An image of the environment "bigworld" was loaded into the program. The width of the image was set to 8.4 cm.

Based on this lines was drawn by hand between the centroids of the rooms navigating any obstacles. This can be seen on figure 6.

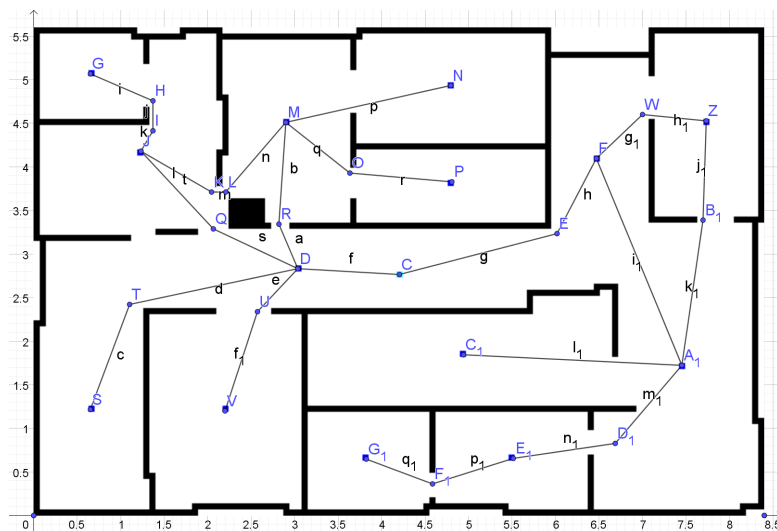


Figure 6: Illustration of paths in geogebra

Geogebra calculates the length of each line. Based on this the total path length from one room to another was found. This can be seen in table gegepgkek

### A.2.2 Test parameters

- World used           bigworld
- Length of world       8.4 cm

### A.2.3 Data

Distribution of marbles			
Distance from room 1 to 2	-2.78	Distance from room 2 to 3	-4.32
Distance from room 2 to 6	-4.58	Distance from room 3 to 4	-3.88
Distance from room 3 to 5	-4.24	Distance from room 3 to 6	-3.46
Distance from room 6 to 7	-6.54	Distance from room 6 to 8	-3.76
Distance from room 7 to 9	-8.02	Distance from room 9 to 10	-2.94
Distance from room 9 to 11	-5.14	Distance from room 10 to 11	-5.64
Distance from room 11 to 12	-5.02	Distance from room 11 to 13	-4.72
Distance from room 13 to 14	-3.56		

Table 4: Distances between rooms



### A.3 The impact of $\epsilon$ on Q-learning performance

The purpose of this test is to show how different values of  $\epsilon$  influences the performance of Q-learning.

#### A.3.1 Description of test

This test was done by performing 100 trials of each selected value of  $\epsilon$  for a range of episodes. The test was conducted using the world "5-room world" seen on figure 7.

The distance punishments are bases on those found in the test in appendix A.2. The distance punishments was scaled with a factor of 1.2 to make the distance a bigger factor in final path.

The probabilities used for each room was found in the test in appendix A.1. These probabilities was divided by the maximal value and scaled by a factor of 20. This was done to ensure that the total reward for entering a room the first time would be positive.

The initial state for all tests was set to room 3 in order to ensure greatest number of possible paths.

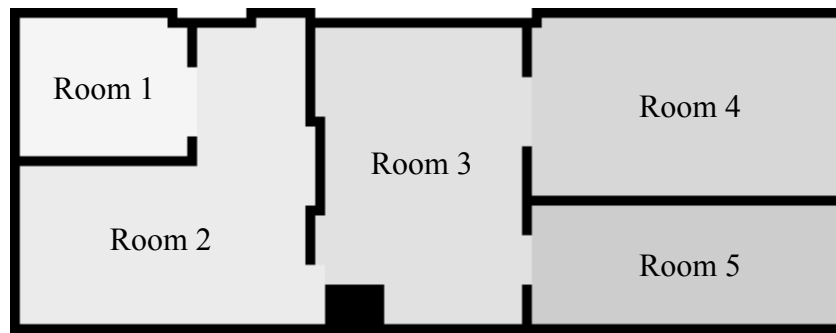


Figure 7: Illustration of "5-room world"

#### A.3.2 Test parameters

In the following tables, the parameters for the test can be seen. In order not to make the agent act completely randomly, it have been chosen to make the test on values of  $\epsilon$  between 0.01 and 0.5.

- World used	5-room world
- Initial room	room 3
- Probabilities based on	50 tests
- Number of tests	100
- Scaling factor distance	1.2
- Scaling factor reward	20
- Learning rate $\alpha$	0.1
- Discount factor $\gamma$	0.9

Tested values of $\epsilon$	
0.01	0.15
0.025	0.2
0.05	0.3
0.075	0.4
0.1	0.5

Table 5: Table of tested the values of  $\epsilon$ . Ranging from 1 % to 50 %

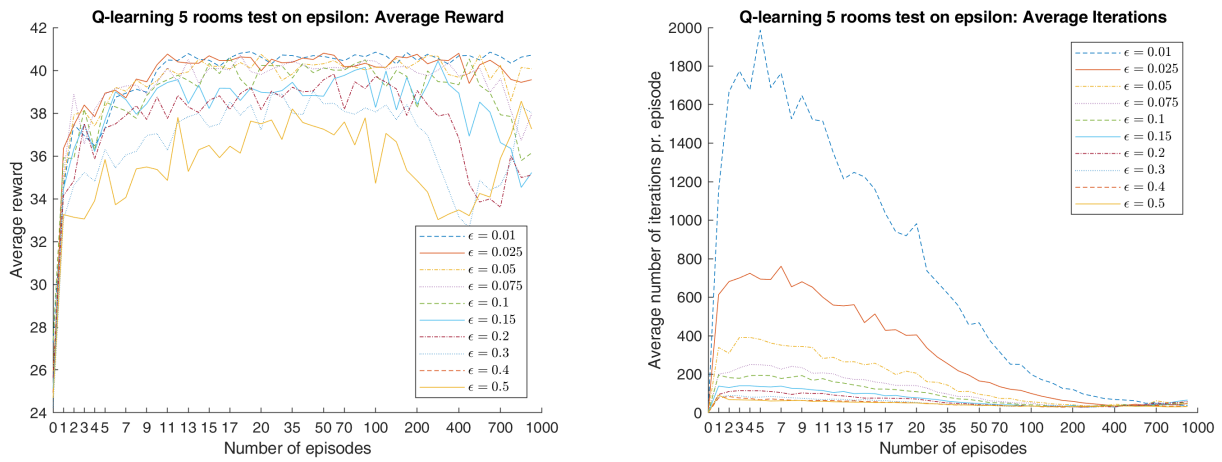
Distance punishments	
Start to room 3	0
Room 1 to room 2	-1.668
Room 2 to room 3	-2.592
Room 3 to room 4	-2.328
Room 3 to room 5	-2.544

Table 6: Table of the distance punishments. The distances are found in the test in appendix A.2

### A.3.3 Data

On figure 8a the average reward for all tests can be seen. It can be seen that the lower the value of  $\epsilon$  the higher average reward, meaning that the policy will converge to the optimal policy.

On figure 8b the average number of iterations per episoden can be seen. It can be seen that the higher the value of  $\epsilon$  the lower the average number of iterations will be. Given the fact that the more random the agent acts, the faster the agent will search alternative paths to the policy, and find its way through the environment.



(a) Plot of how the average reward develops as a function of the number of episodes for each value of  $\epsilon$

(b) Plot of how the average number of episodes develops as a function of the number of episodes for each value of  $\epsilon$

Figure 8: Plots of both the average reward and average number of iterations pr. episode for each value of  $\epsilon$

A value of 0.05 have been chosen as the best compromise between a high average reward and a low average number of iterations per episode.

### A.3.4 Conclusion

It can be concluded that the smaller the value of  $\epsilon$  the higher the average reward will be for a given number of episodes.

It can as well be concluded that the higher the value of  $\epsilon$  the lower the number of iterations per episode will be.

The value of 0.05 was chosen as the best compromise.

## A.4 The impact of $\alpha$ on Q-learning performance

The purpose of this test is to show how different values of  $\alpha$  influences the performance of Q-learning.

### A.4.1 Description of test

This test was done by performing 100 trials of each selected value of  $\alpha$  for a range of episodes. The test was conducted using the world "5-room world" seen on figure 9.

The distance punishments are bases on those found in the test in appendix A.2. The distance punishments was scaled with a factor of 1.2 to make the distance a bigger factor in final path.

The probabilities used for each room was found in the test in appendix A.1. These probabilities was divided by the maximal value and scaled by a factor of 20. This was done to ensure that the total reward for entering a room the first time would be positive.

The initial state for all tests was set to room 3 in order to ensure greatest number of possible paths.

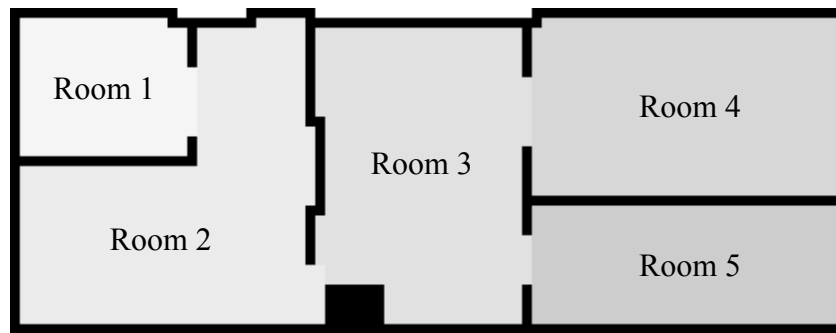


Figure 9: Illustration of "5-room world"

### A.4.2 Test parameters

In the following tables, the parameters for the test can be seen. It have been chosen to make the test on values of  $\alpha$  between 0.01 and 0.5.

- World used	5-room world
- Initial room	room 3
- Probabilities based on	50 tests
- Number of tests	100
- Scaling factor distance	1.2
- Scaling factor reward	20
- Randomness factor $\epsilon$	0.05
- Discount factor $\gamma$	0.9

Tested values of $\alpha$	
0.01	0.15
0.025	0.2
0.05	0.3
0.075	0.4
0.1	0.5

Table 7: Table of tested the values of  $\alpha$ . Ranging from 0.01 to 0.5

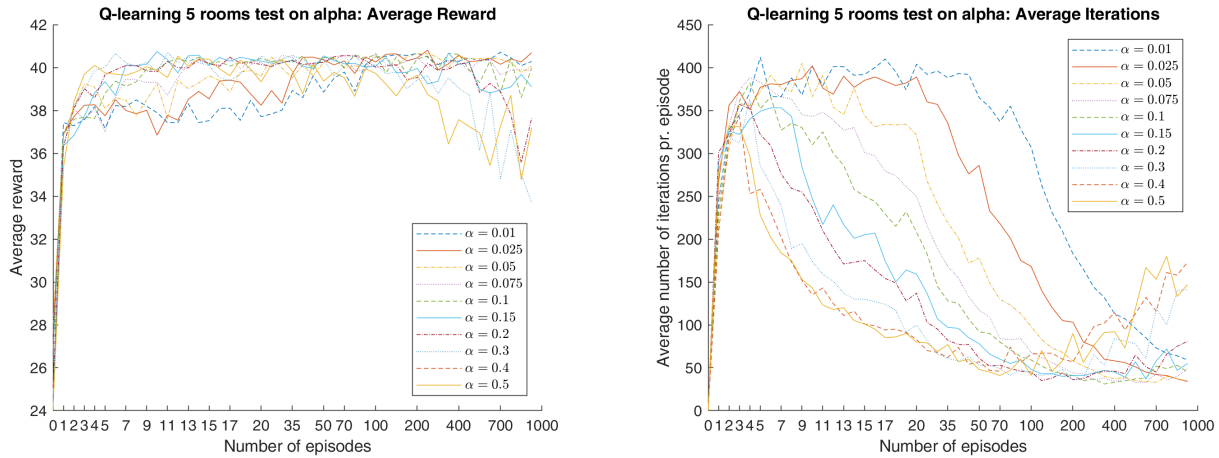
Distance punishments	
Start to room 3	0
Room 1 to room 2	-1.668
Room 2 to room 3	-2.592
Room 3 to room 4	-2.328
Room 3 to room 5	-2.544

Table 8: Table of the distance punishments. The distances are found in the test in appendix A.2

### A.4.3 Data

On figure 10a the average reward for all tests can be seen. It can be seen that the lower the value of  $\alpha$  the higher average reward, meaning that the policy will converge to the optimal policy. Whereas a higher value of  $\alpha$  will result in a lower average reward. It would also seem like the average reward are oscillating more.

On figure 8b the average number of iterations per episoden can be seen. It can be seen that the higher the value of  $\alpha$  the lower the average number of iterations will be. This is due to the fact, that the lower the learning rate are, the slower the algorithm learns, and will therefore need more steps to get to the goal.



(a) Plot of how the average reward develops as a function of the number of episodes for each value of  $\epsilon$

(b) Plot of how the average number of episodes develops as a function of the number of episodes for each value of  $\alpha$

Figure 10: Plots of both the average reward and average number of iterations pr. episode for each value of  $\alpha$

A value of 0.025 have been chosen as the best compromise between a high average reward and a low average number of iterations per episode.

### A.4.4 Conclusion

It can be concluded that the smaller the value of  $\alpha$  the higher the average reward will be for a given number of episodes. It can as well be concluded that the higher the value of  $\alpha$  the lower the number of iterations per episode will be. The value of 0.025 was chosen as the best compromise.

## A.5 Motion planners in action

The purpose of this test is to find out how well the motion planners work in an actual environment with obstacles. The goal is to lead the robot from an initial position to a target location.

### A.5.1 Description of test

The initial position of the robot will be in origo of the environment of the gazebo simulator. Here 14 test of both the tangent bug algorithm and model based planner will be conducted for the 14 rooms of the *bigworld* map. The idea is to test the success rate of finding rooms, the robots distance to the closest obstacle on the path to the goal and the distance travelled along the way. The model based planer will be tested with a higher speed than the tangent bug algorithm because it will be further away from obstacle most of the time and thereby less likely of hitting obstacles.

### A.5.2 Test parameters

- World used                      bigworld
- Speed of tangent bug        -1 to 1
- Speed of model planer      -2 to 2
- Number of tests                14

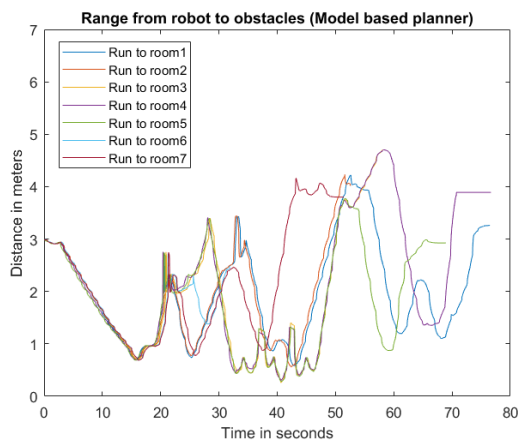
### A.5.3 Data

<b>Modelbased</b>					
	Successrate (%)	Distance traveled (m)	Time (s)	Approx velocity (m/s)	Distance to obstacle (m)
Room 1	100	63.087	76.474	0.825	1.990
Room 2	100	43.331	52.589	0.824	1.876
Room 3	100	44.450	58.178	0.764	1.886
Room 4	100	57.150	76.660	0.745	2.157
Room 5	100	50.800	68.931	0.737	1.809
Room 6	100	23.576	31.978	0.737	1.822
Room 7	100	41.920	51.427	0.815	2.127
Room 8	100	25.400	33.170	0.766	2.401
Room 9	100	31.044	36.554	0.849	2.752
Room 10	100	44.450	52.626	0.845	2.435
Room 11	100	37.355	35.697	1.046	2.628
Room 12	100	59.851	69.724	0.858	2.966
Room 13	100	57.029	67.993	0.839	3.094
Room 14	100	76.906	87.035	0.884	2.893
<b>Average</b>	<b>100</b>	<b>46.882</b>	<b>57.074</b>	<b>0.824</b>	<b>2.345</b>

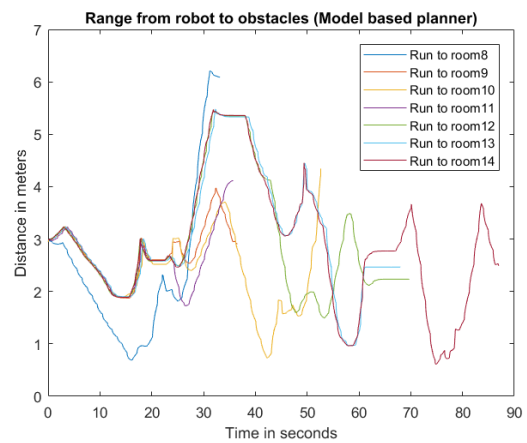
Table 9: Table over test data for the model based planner

Sensorbased					
	Successrate (%)	Distance traveled (m)	Time (s)	Approx velocity (m/s)	Distance to obstacle (m)
Room 1	0	-	-	-	-
Room 2	100	31.044	87.956	0.353	1.613
Room 3	100	23.283	79.409	0.293	1.683
Room 4	100	38.806	135.232	0.287	1.316
Room 5	100	28.222	415.227	0.068	1.033
Room 6	100	8.467	28.323	0.299	3.050
Room 7	0	-	-	-	-
Room 8	100	23.283	79.528	0.293	1.724
Room 9	100	23.989	83.720	0.287	1.193
Room 10	100	38.806	123.614	0.314	1.432
Room 11	100	35.278	102.771	0.343	1.320
Room 12	0	-	-	-	-
Room 13	0	-	-	-	-
Room 14	0	-	-	-	-
<b>Average</b>	<b>64.286</b>	<b>27.909</b>	<b>126.198</b>	<b>0.282</b>	<b>1.596</b>

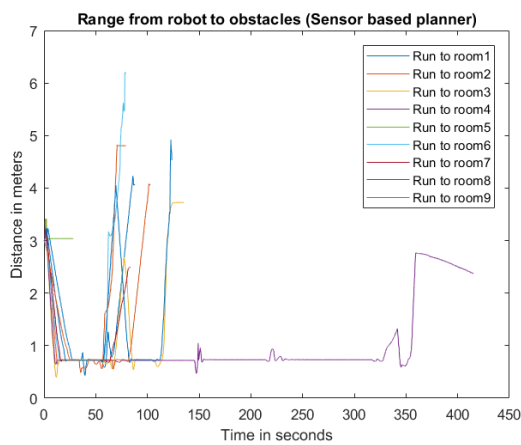
Table 10: Table over test data for the sensorbased planner



(a) Illustration of the distance to the closest obstacle for room 1-7 for the model based planner

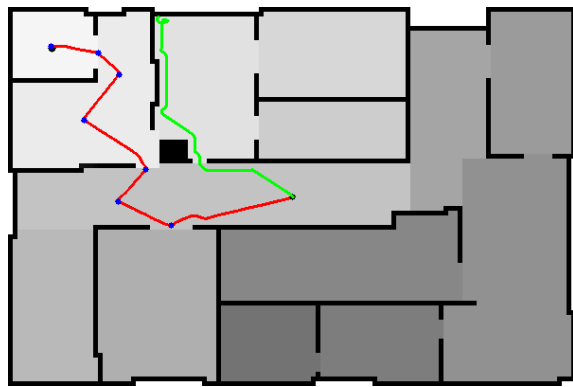


(b) Illustration of the distance to the closest obstacle for room 8-14 for the model based planner

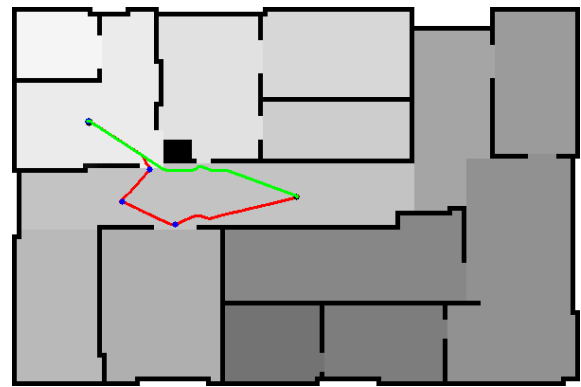


(c) Illustration of the distance to the closest obstacle for room 1-14 for the sensor based planner

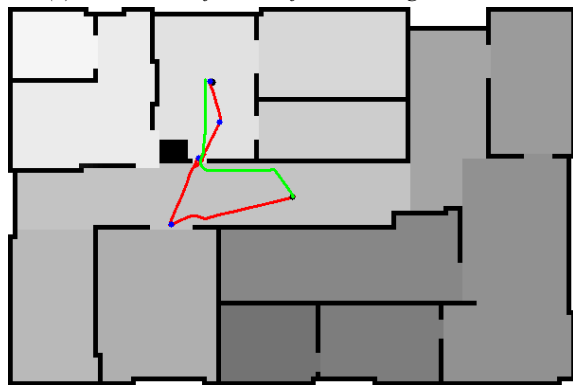
Figure 11: Illustration of the distance to the closest obstacle for room 1-14 for both planners



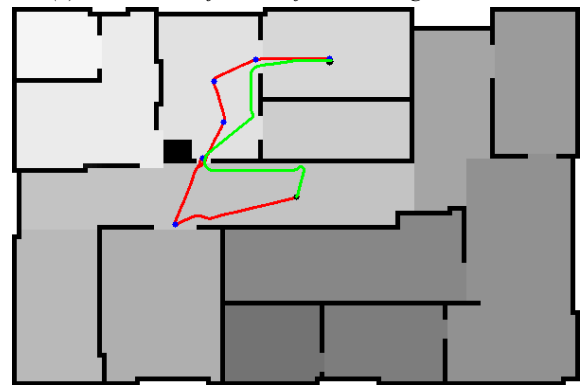
(a) Illustration of the run from the origin to room 1



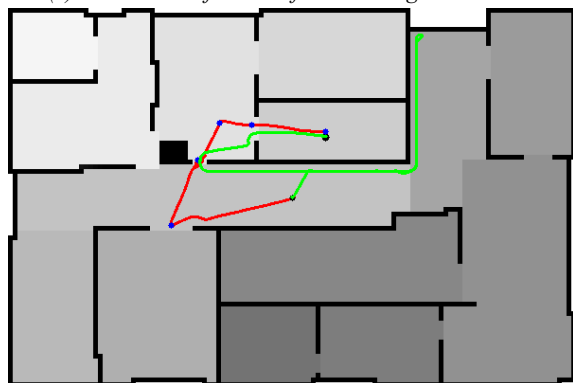
(b) Illustration of the run from the origin to room 2



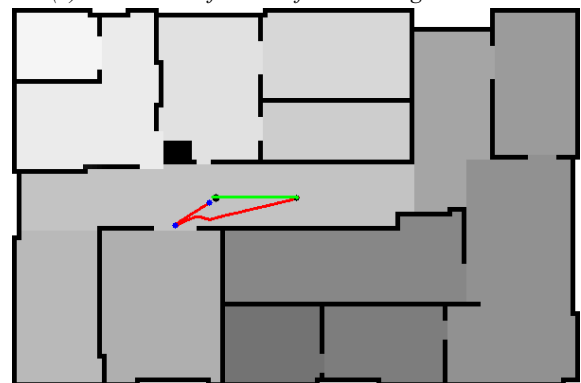
(c) Illustration of the run from the origin to room 3



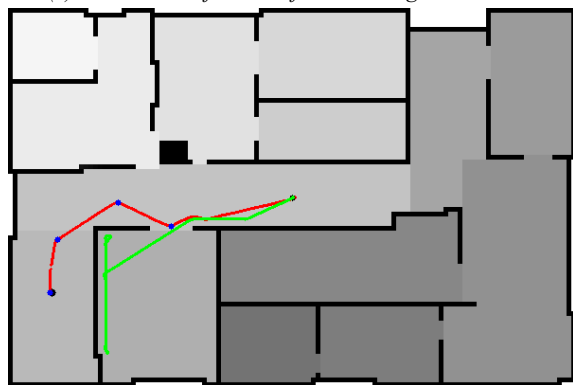
(d) Illustration of the run from the origin to room 4



(e) Illustration of the run from the origin to room 5



(f) Illustration of the run from the origin to room 6

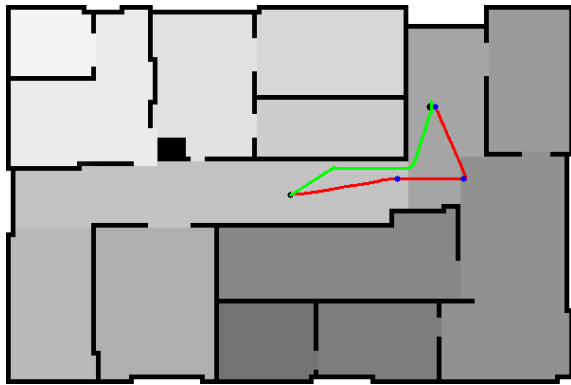


(g) Illustration of the run from the origin to room 7



(h) Illustration of the run from the origin to room 8

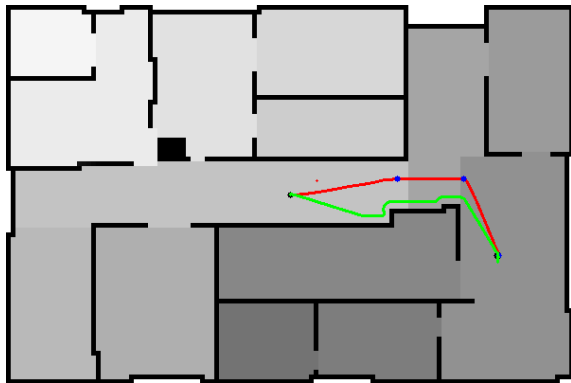
Figure 12: Illustration of the run from the origin to room 1-8 for both the model based (red line) and the sensor based (green line)



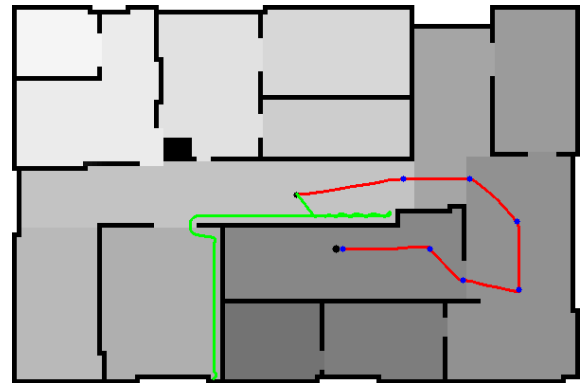
(a) Illustration of the run from the origin to room 9



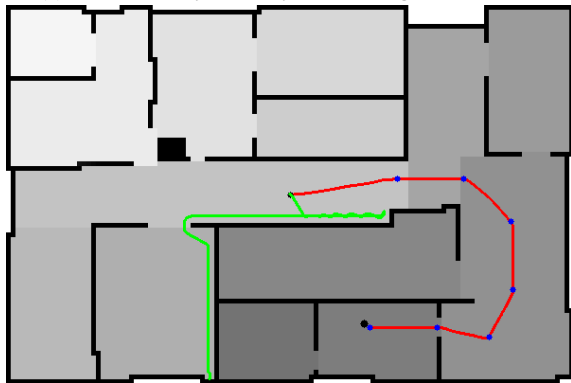
(b) Illustration of the run from the origin to room 10



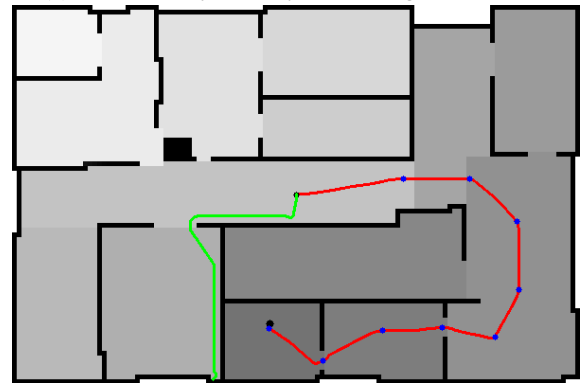
(c) Illustration of the run from the origin to room 11



(d) Illustration of the run from the origin to room 12



(e) Illustration of the run from the origin to room 13



(f) Illustration of the run from the origin to room 14

Figure 13: Illustration of the run from the origin to room 9-14 for both the model based (red line) and the sensor based (green line)



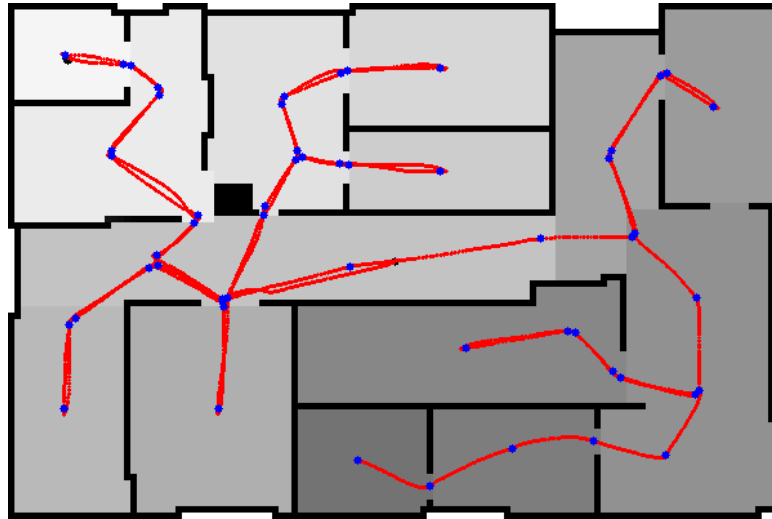


Figure 14: Illustration of the run visiting all rooms starting from the origin for the model based planner

#### A.5.4 Conclusion