

# COMP0037 2021/22 Coursework 1

COMP0037 Teaching Team

18th January 2022

## Overview

- Coursework 01 Release Date: Tuesday, 18th January 2022; Latest Revision Date: 17th January 2022
- **Assignment Due Date: Wednesday, 16th February 12:00 (Midday UK time) 2022**
- Weighting: 40% of the module total
- **Final Submission Format.** Each group must submit *three* things:
  1. A zip file (which contains the source code implemented to tackle this coursework).  
The name of the zip file must be of the form `COMP0037_CW1_GROUP_g.zip`, where `g` is the letter code of your group.
  2. A report in PDF format. It will be named `COMP0037_CW1_GROUP_g.pdf`.
  3. A video which describes your answers. You may use a variety of video formats, but the video must be playable by [VLC](#). It will be named `COMP0037_CW1_GROUP_g.[ext]`.

The total marks for each question are written in the form [*X* marks]. Approximately 67% of the total coursework marks are from your analysis and 33% from your coding of the algorithms. Therefore, please take time and care with your writing and analysis of solutions.

For the submitted code, you must use the notation variable names and code provided in the lectures. If you do not use these, we will consider the work a potential plagiarism case and investigate further.

You will need to implement additional routines to support the functionality required in these questions. The code to implement this is available on Moodle. It is based upon but extends and refactors the code provided in the labs. The general areas requiring adjustment will be shown using comments.

We expect you to take care in writing your report. For example, figures and graphs need to have captions and be numbered and labelled. Equations captured by a screenshot, from the slides or elsewhere, and pasted into the document will not be accepted. When we ask you to write some code, provide an explanation in your report of what you wrote.

The final mark will be computed by summing all the marks awarded on individual questions together and dividing by the mark total.

Neither the video nor the code will be independently marked, but both must be provided before the submission deadline. Both will be checked.

## Use Case

All the questions in this coursework are based on an example of a robot that automatically cleans an airport arrivals area. Airports are generally very busy places and require constant cleaning.

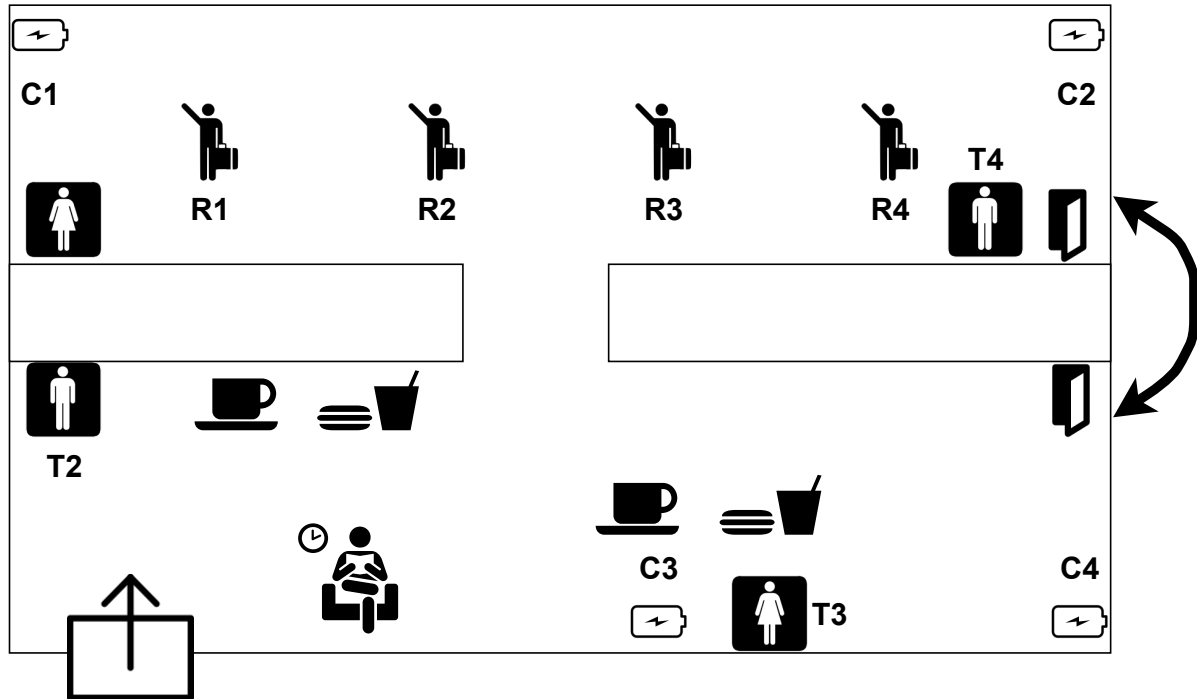


Figure 1: Overview of the arrivals area scenario. This is very loosely based on Heathrow Terminal 4. See text for description. Note this figure is not drawn to scale. The actual scenario is generated by the function `full_scenario` in the provided code (script: `scenarios`).

Figure 1 shows an overview of the environment the robot operates in. Passengers arrive at the top of the map, collect their luggage from the baggage reclaim areas (R1–R4) and pass through customs. Beyond there are various food vendors and a waiting area. Passing through the customs area is a high-cost action because of the disruption the robot would cause. Instead, the robot uses a secret door on the right-hand side of the map that only it can access. Because this door needs to be securely opened and closed, the travel costs associated with travelling through this door are much higher than driving in open space.

There are four toilet blocks (T1–T4). T1 and T4 are before customs, and T2 and T3 are after customs. The cleaning robot is powered by a battery which becomes depleted over time. Four charging stations are located throughout the terminal (C1–C4). Each charging location has a different charging rate resulting in a different charging time.

## Questions

1. Although four outlets are provided for the robot to charge its battery, the different chargers have different charging rates. The charging rates are modelled as Gaussian distributions with variance 1. The true charging rates are shown in Table 1.

The robot has been tasked with learning the best charging location. Still, the only information it can receive is visiting charging stations and measuring the charging rate it observes.

Charging Location	Mean Charging Rates (Amps)
1	4
2	4.1
3	3.9
4	4.2

Table 1: Charging information for locations.

- a. Briefly describe how this problem can be expressed as a bandit problem and specify the number of arms. Describe what the epistemic and aleatoric uncertainties are for this problem.

[6 marks]

- b. Implement a multi-arm bandit simulation to simulate the charging rates by the different stations. Describe your code. Generate violin plots to visualise the locations using the data generated from your simulator.

[6 marks]

- c. Define what regret is and present the formula used to calculate it. Explain why it is used to assess the performance of bandit problems. Although we can use regret for this problem, why might it not be an available tool for other problems?

[5 marks]

- d. Describe the  $\epsilon$ -greedy strategy. Prove that, if there are  $C$  charging stations the robot could visit, the *long-term average probability* of visiting the optimal station is

$$1 + \left( \frac{1}{C} - 1 \right) \epsilon.$$

For very short runs, the number of correctly pulled handles could be significantly lower than this value. Propose an explanation for why this is the case and show how it could be modelled using the equation above.

[15 marks]

- e. Extend your implementation of the bandit problem to support Monte Carlo analysis (in which you can run the whole experiment with the agent multiple times and average the rewards across all those runs).

By varying  $\epsilon$  and experimenting with the number of Monte Carlo runs, demonstrate that the equation above appears to be correct.

[15 marks]

- f. A number of ways to improve the performance of  $\epsilon$ -greedy has been proposed. Describe one approach, implement it, and demonstrate its performance using graphs of both cumulative regret and percentage of optimal actions taken.

[15 marks]

- g. Briefly describe Upper-Confidence-Bound (UCB) Action Selection, including the intuition of why this algorithm was proposed, and outline the equations used to implement it.

[5 marks]

- h. Investigate the performance of the UCB for several different choices of the degree of exploration. What trend do you observe, and what do you think is the value which gives the best performance. How do you think the degree of exploration relates to the variance of the charging rates?

[10 marks]

[Question total 77 marks]

2. a. Describe the *breadth-first search* and *depth-first search* algorithms. What are the advantages and disadvantages of each?

[4 marks]

The script `compare_path_planning_algorithms.py` causes the robot to plan a path between each rubbish bin in turn in the environment.

- b. Evaluate the performance of the breadth-first and depth-first planners for this case. Your analysis should include quantitative measures (e.g., the sum of all path costs and the sum of all cells visited when planning the paths) and a qualitative discussion of any unusual behaviour observed. Use screenshots to illustrate unusual behaviours and explain why they arise.

[10 marks]

- c. Describe *Dijkstra's algorithm* and explain the two key changes it has compared with breadth-first search. What is the purpose of each of these modifications?

[6 marks]

- d. Modify the class `DijkstraPlanner` to implement Dijkstra's algorithm and explain your implementation. Compare the results with breadth-first search with regards to path cost and the number of cells visited.

[18 marks]

Suppose a path takes a cell  $i$  into an adjacent cell  $j$ . The current  $l$ -stage additive cost is modelled as the Euclidean distance between the cells,

$$l(s_i, s_j) = |s_j - s_i|.$$

However, this does not take account of traversability constraints. These are modelled by

$$l(s_i, s_j) = \alpha(s_j) |s_j - s_i|,$$

where  $\alpha(s_j)$  is a penalty factor based on the type of  $j$ 's cell. Its value is given by

$$\alpha(s_j) = \begin{cases} 5 & \text{if } s_b \text{ is part of the secret door} \\ 100 & \text{if } s_b \text{ is part of the customs area} \\ 1 & \text{otherwise} \end{cases}$$

- e. Modify the implementation of `compute_transition_cost` in `airport_map` to take account of  $\alpha(s_j)$ . Modify the code to implement this change and explain what changes you made. Evaluate how this changes the performance of Dijkstra's algorithm.

[9 marks]

- f. Describe the key ideas behind the A\* algorithm. What does it mean for a heuristic to be admissible, and why is it important? Explain why the Euclidean distance is an admissible heuristic, but the squared Euclidean distance is not.

[8 marks]

- g. What do you think is the *optimal heuristic* that A\* could use? Could it be used in practice? Explain your answer.

[5 marks]

- h. Implement the A\* algorithm in the `AStarPlanner` class with the Euclidean heuristic and describe your implementation. Compare results with the earlier algorithms with regards to path cost and the number of cells visited.

[12 marks]

The A\* planning algorithm with an admissible heuristic will be used for all future robot operations.

The designers need to design a recharge policy for the robot: when its battery has become sufficiently low, the robot needs to drive to a recharging station and recharge itself. This can be expressed in the policy

$$a_t = \pi(s_t),$$

where  $s_t$  is the state where the robot decides it needs to recharge and  $a_t \in [1, 2, 3, 4]$  is the charging station the robot will go to. Suppose the  $i$ th charging station is at  $s_i$  and the recharge is Gaussian distributed with mean  $\mu_i$  and covariance  $\sigma_i^2$ . The action function is

$$q_\pi(s, a) = \mu_a - l(s, s_a).$$

The values for  $\mu_i$  and  $\sigma_i^2$  are known. The quantities have been rescaled to match the units with the total path length. Note that the means of the charging stations differ from those in Q1.

- i. Propose an algorithm to compute  $\pi_*(s)$ .

[6 marks]

*Hint:* You do not need to use policy or value iteration.

- j. Implement your algorithm in the `compute_recharging_policy` script. The resulting policy must be stored in an instance of the `ChargingPolicy` class. Explain your implementation.

Display the learned policy using the `ChargingPolicyDrawer`. Discuss the results. Are they as you expect?

[12 marks]

- k. Pick an inadmissible heuristic for your A\* algorithm. Explain the heuristic used. Compute the new policy and display the results. Comment on whether the use of the inadmissible heuristic made a significant difference or not.

[4 marks]



1. Suppose the means and covariances for each charging station  $(\mu_i, \sigma_i^2)$  were not known. How would you modify your algorithm to handle this case? You do not have to implement the resulting algorithm.

[6 marks]

[100 marks]

The engineers realise that the robot is actually very poorly built and very poorly calibrated. As a result, when commanded to move in a given direction, it is not guaranteed to move in that direction. Rather, as illustrated in Figure 2, the robot moves in the nominal direction with a probability  $p$  or one square on either side of the nominal direction with a probability  $q = 0.5(1 - p)$ .

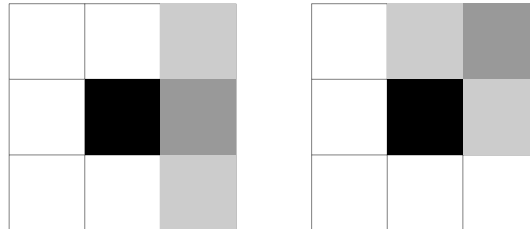


Figure 2: Examples of the errors in the process model. In both cases, the original robot position is the black square. The nominal direction of motion is the dark grey square. The errors in the robot mean that the robot will either end up in the dark grey square or either light grey square.

The traversability costs are the same as in question 2. Because the robot's motion is uncertain, it might collide with objects in the environment. If the robot collides with the luggage reclaim area, it experiences a cost of -10. If it collides with any other object, it incurs a cost of -1. The robot does not need to be recharged for the duration of the task of interest.

3. The engineers determine that this system can be modelled as a Finite Markov Decision Process (FMDP).

a. Define what a Finite Markov Decision Process (FMDP) is, and describe its main components.

[5 marks]

b. Consider the case where the robot starts at a known start location and has to reach a known goal location. When the robot reaches the goal, the episode ends. Explain how this can be modelled as an FMDP. Your answer should relate the problem to each component of an FMDP you identified in part (a).

[15 marks]

The class `PolicyIterator` will implement the policy iteration algorithm. The provided code only implements the policy evaluation step.

- c. Use the code in `q3_cef.py` to apply the policy iteration algorithm to the system. The policy evaluation algorithm using the initial policy and  $\gamma = 1$  will not converge. Why is this the case? Suggest a way to fix this and show that the change causes policy evaluation to converge.

[7 marks]

- d. Modify the class `PolicyIterator` to complete the implementation of the policy improvement step. Describe your implementation.

[12 marks]

- e. Run the script `q3_cef.py`. Compute the policies and value functions for  $p = 1$ ,  $p = 0.9$ ,  $p = 0.6$  and  $p = 0.3$ . In each case, plot the policy and the value function for each cell. Discuss the behaviour of the policy and value function for the different values of  $p$ .

[12 marks]

For the rest of this coursework,  $p = 0.8$ .

- f. Modify the `PolicyIteration` class to record the number of iterations required internally for each policy evaluation and the total number of iterations required for a policy to converge. Explain your code. Run the script `q3_cef.py`, present the results, and discuss any trends or behaviours you see.

[12 marks]

- g. Run the script `q3_g.py` to apply your algorithm to the full problem. Present the generated policy, value function, and policy iteration algorithm statistics.

By referring to the policy and the airport map, discuss the uncertainty's impact on the computed policy.

[8 marks]

- h. Complete the implementation of the value iteration in the `ValueIterator` class and describe your code.

[12 marks]

- i. Modify the script `q3_i.py` to compare policy iteration and value iteration on the full map and discuss the difference between the two algorithms. This includes the number of iterations, required and differences in generated policy. Which algorithm would you recommend and why?

[10 marks]

[Total 94 marks]

The actual robot will be configured to operate sequentially, moving between rubbish bins and toilets in a pre-defined sequence. As the robot moves through the sequence, the engineers determine that its battery will become depleted. If the battery becomes flat, the robot will have to be rescued, which would incur a large negative cost.

After it has visited one cleaning site, the engineers propose that the robot decide to either move directly to the following cleaning site or visit a recharge station, recharge itself, and then move to the following cleaning site. The recharge stations provide Gaussian-distributed rewards, similar to Q1.

4. Describe, but do not implement, how you would represent this as an FMDP and how policy or value iteration will be applied. Your description should be in reference to the components of an FMDP.

[12 marks]

[Total 12 marks]

[Coursework Total 283 marks]