

Q-Learning for World Grid Navigation

EE5904 / ME5404 Project II

TA: Zhaoshan Liu
e0575844@u.nus.edu

Contents



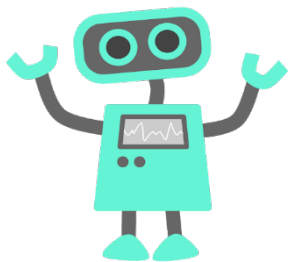
| | |
|-----------------------|-------|
| ○ Project Description | 00:38 |
| ○ Implementation | 02:31 |
| ○ Task 1 | 05:30 |
| ○ Task 2 | 06:54 |
| ○ Assessment | 07:12 |
| ○ Submission | 07:42 |
| Appendix | |
| ○ MATLAB Course | 08:24 |
| ○ Recap | 08:36 |

Project Description

- Task
- State Transition
- Reward Function
- Learning

Task

Using **Q-learning** with ϵ -greedy exploration. The robot is to move from the initial state ($s = 1$) to the goal state ($s = 100$) with the maximum total reward of the trip.



| | | | | | | | | | |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| START | | | | | | | | | |
| 1 | 11 | ... | ... | ... | ... | ... | ... | ... | 91 |
| 2 | 12 | ... | ... | ... | ... | ... | ... | ... | 92 |
| 3 | 13 | | | | | | | . | 93 |
| 4 | . | | | | | | | . | 94 |
| 5 | . | | | | | | | . | 95 |
| 6 | . | | | | | | | . | 96 |
| 7 | . | | | | | | | . | 97 |
| 8 | . | | | | | | | . | 98 |
| 9 | . | | | | | | | 89 | 99 |
| 10 | ... | ... | ... | ... | ... | ... | ... | 90 | 100 |

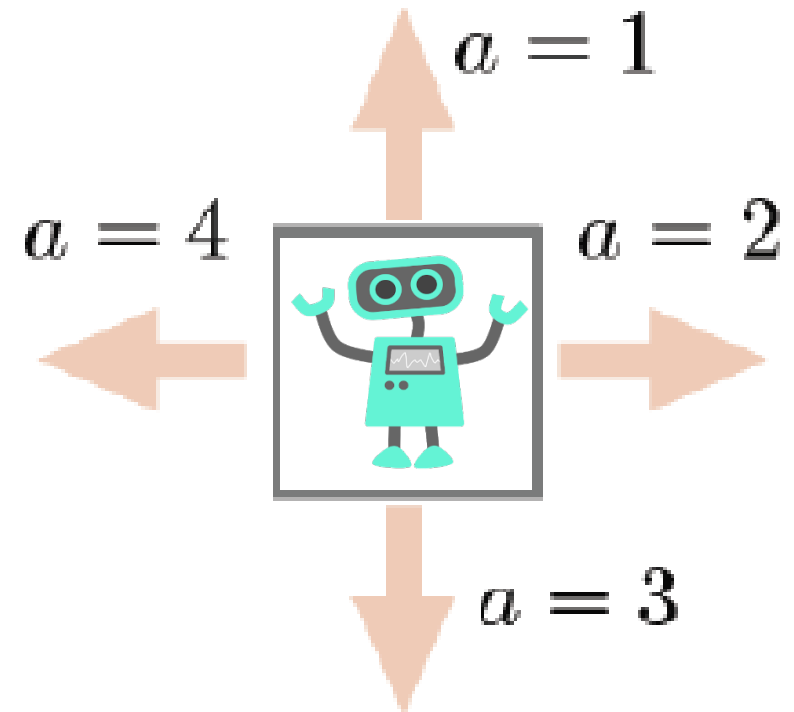


State Transition

At a state, the robot can take one of **four** actions

Deterministic Model

Use dynamic programming methods to find the optimal policy



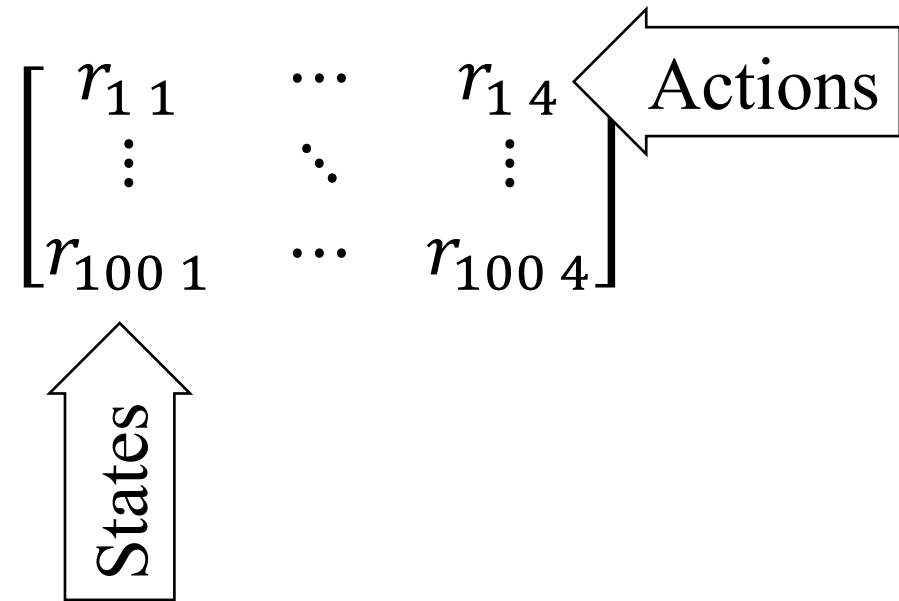
Reward Function

Files

- Task 1 → “reward” in “**task1.mat**”
- Task 2 → “qevalreward” in “**qeval.mat**”

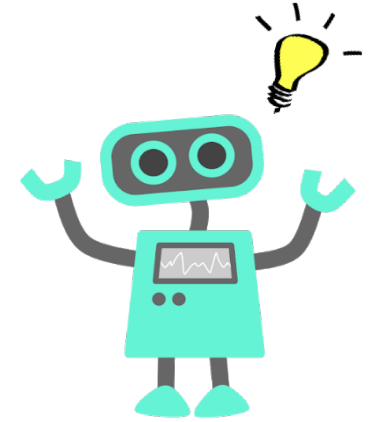
Reward Matrix:

- Column → Action
- Row → State
- 100 x 4



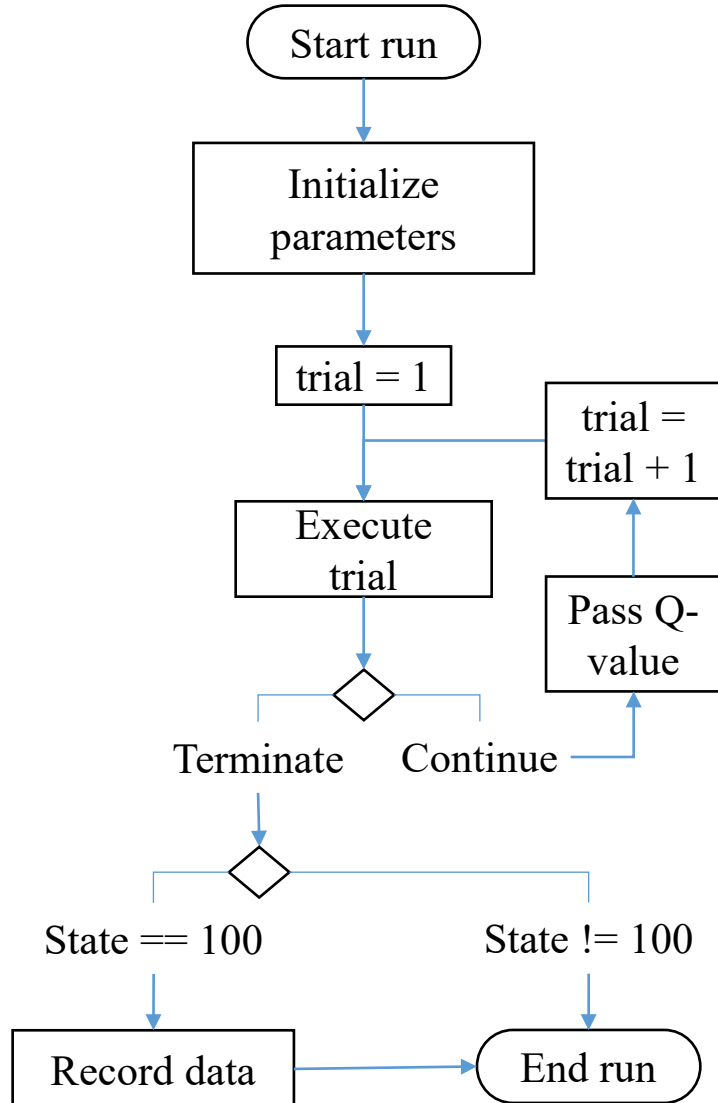
Learning

- The robot learns in 1 run
- One run consists of the N trials
- Each run starts with a set of initial values of the Q-function (100 x 4 matrix)
- Each trial starts when the robot moves from state 1
- Each trial ends when the robot reaches state 100
- The Q values are passed to the next trial
- Each run ends when the Q values converge to the optimal values

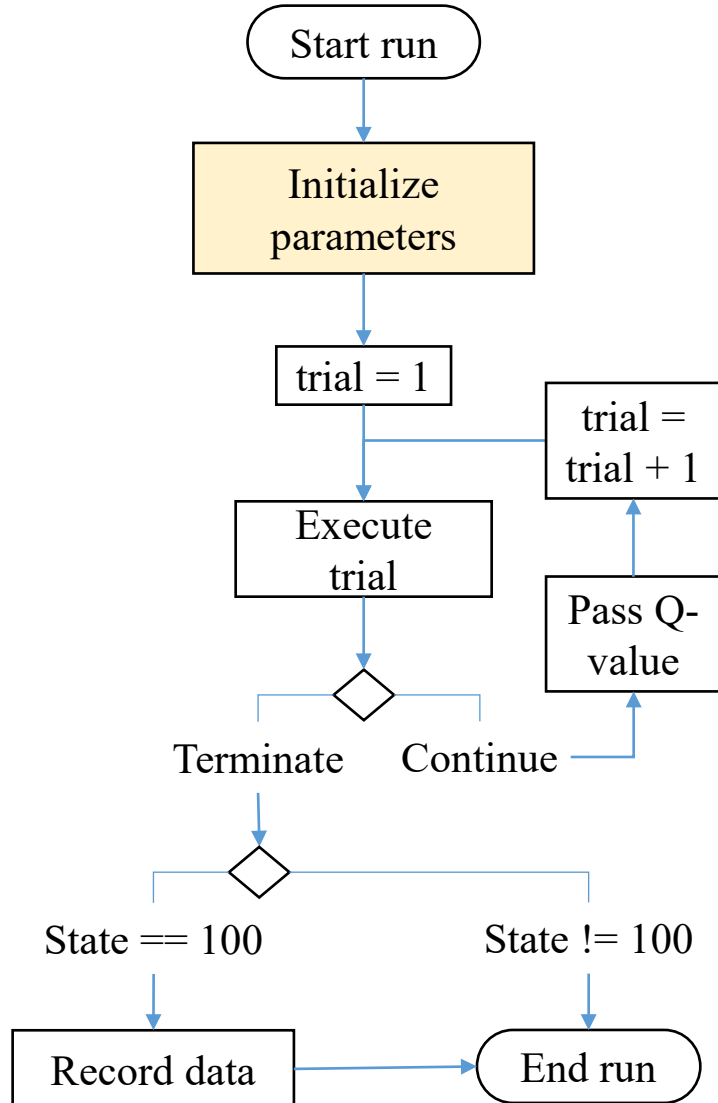


Implementation

Implementation



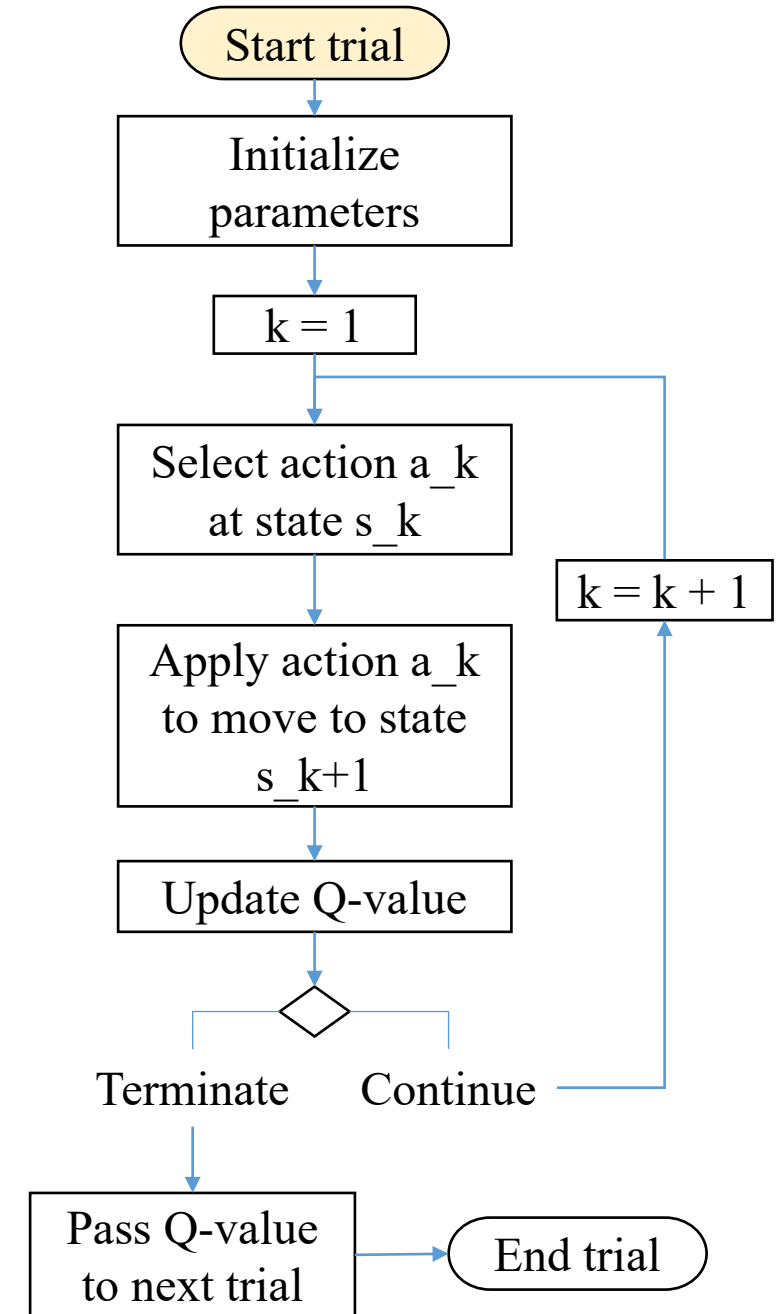
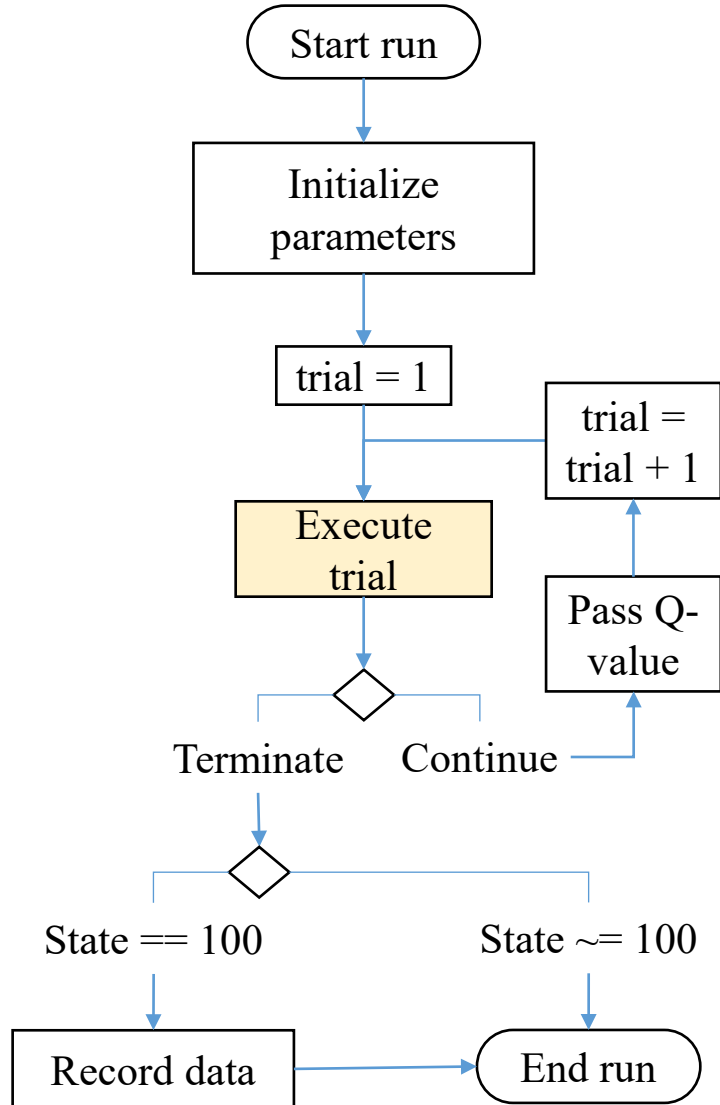
Implementation



Parameters:

- Initial Q-function $Q_1 \leftarrow 0$ >> Optional
- Trial = 1
- threshold of error of Q-values between trials

Implementation



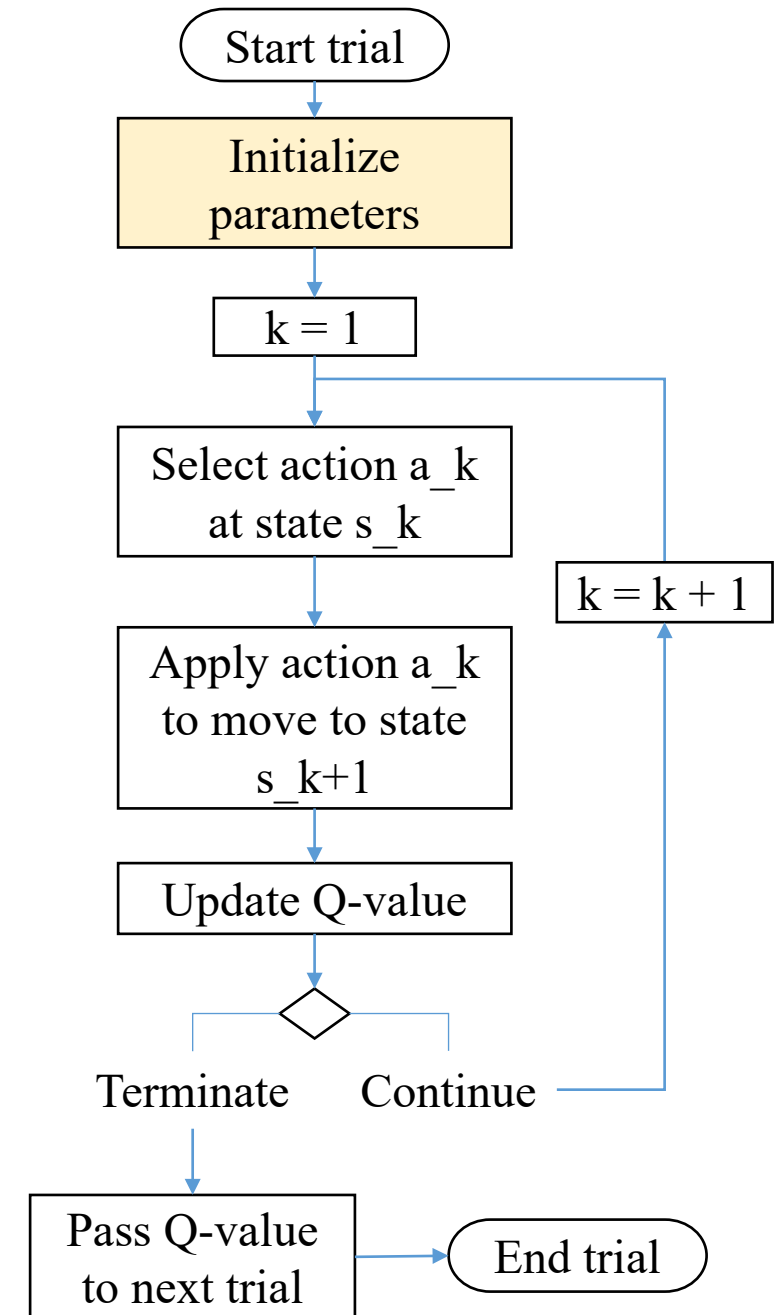
Implementation

Parameters:

- Discount factor γ
- Exploration probability ϵ_k
- Initial Q-function Q_1

From previous trial, if any

- Learning rate $\alpha_k = \epsilon_k$
- Initial state $s_1 = 1$
- Time step $k = 1$



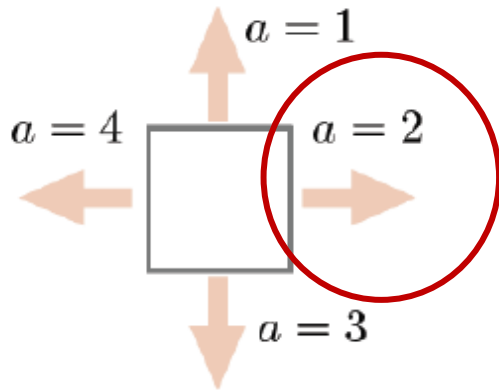
Implementation

Example:

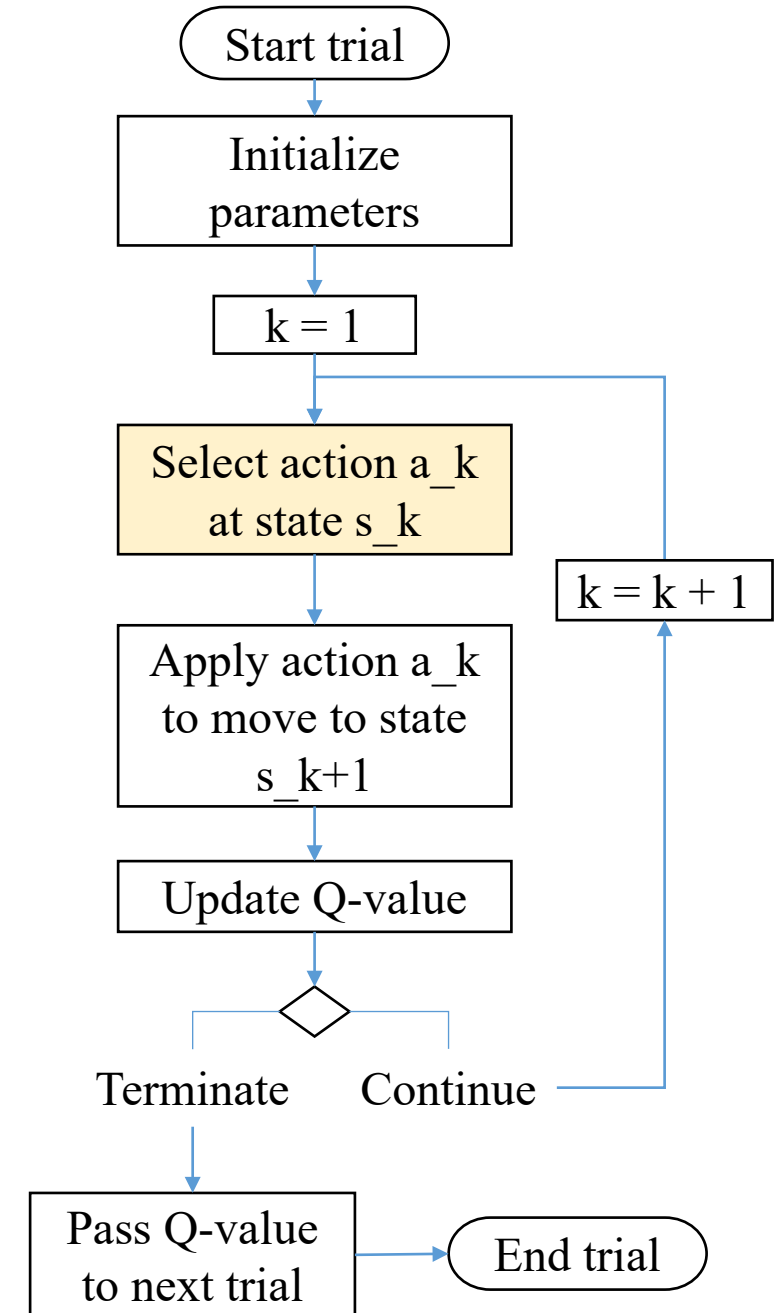
For $k = 3$,

- $\epsilon = 1 - \frac{1}{k}$
- $s_3 = 11$
- Current best action is 2
- Exploitation: $a_3 = 2$ each has $1 - \epsilon = \frac{1}{3}$ probability to be selected
- Exploration: $a_3 = 3, 4$ each has $\epsilon = \frac{2}{(2)3}$ probability to be selected
- $a_3 = 1$ cannot be selected due to the boundary

$$a_k = \begin{cases} a \in \arg \max_{\hat{a}} Q_k(s_k, \hat{a}) & \text{with probability } 1 - \epsilon_k \\ \text{an action uniformly randomly} \\ \text{selected from all} & \text{with probability } \epsilon_k \\ \text{other actions} & \\ \text{available at state } s_k & \end{cases}$$



| | | |
|---|----|-----|
| 1 | 11 | ... |
| 2 | 12 | ... |
| 3 | 13 | |
| 4 | . | |

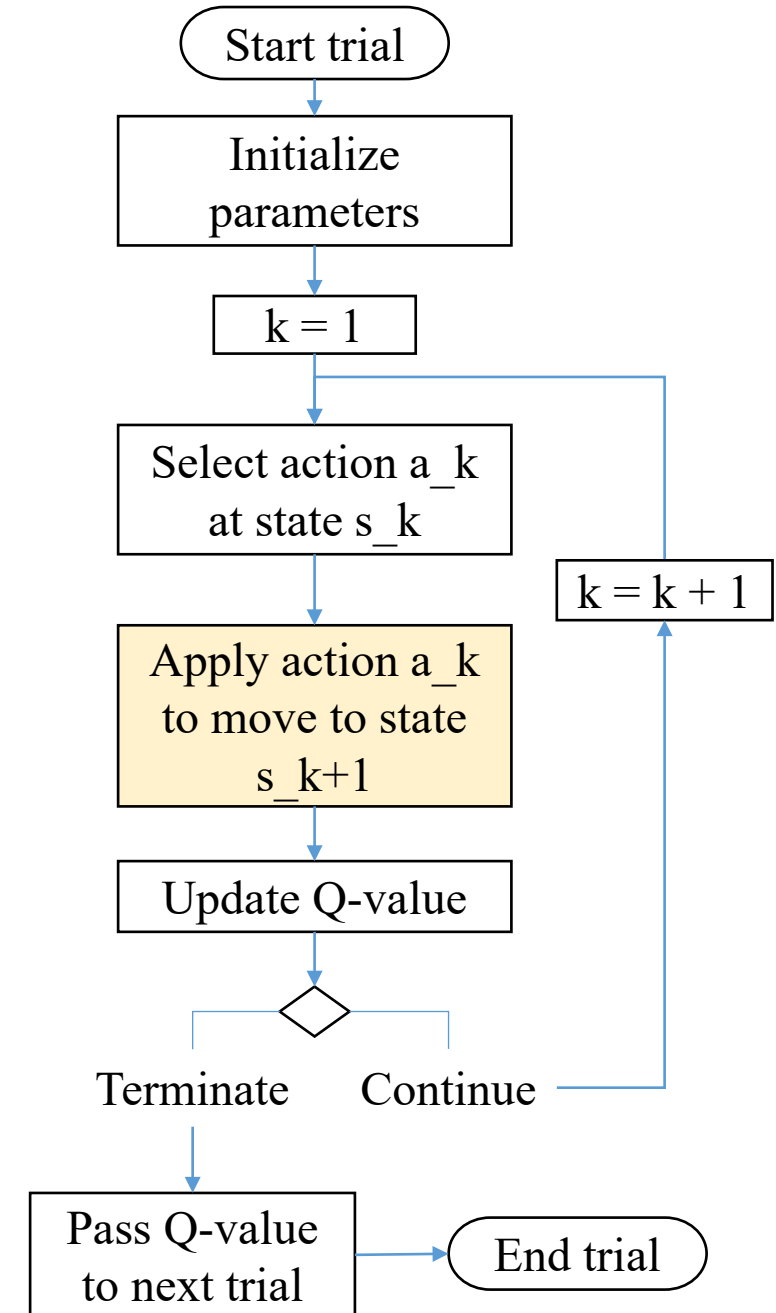
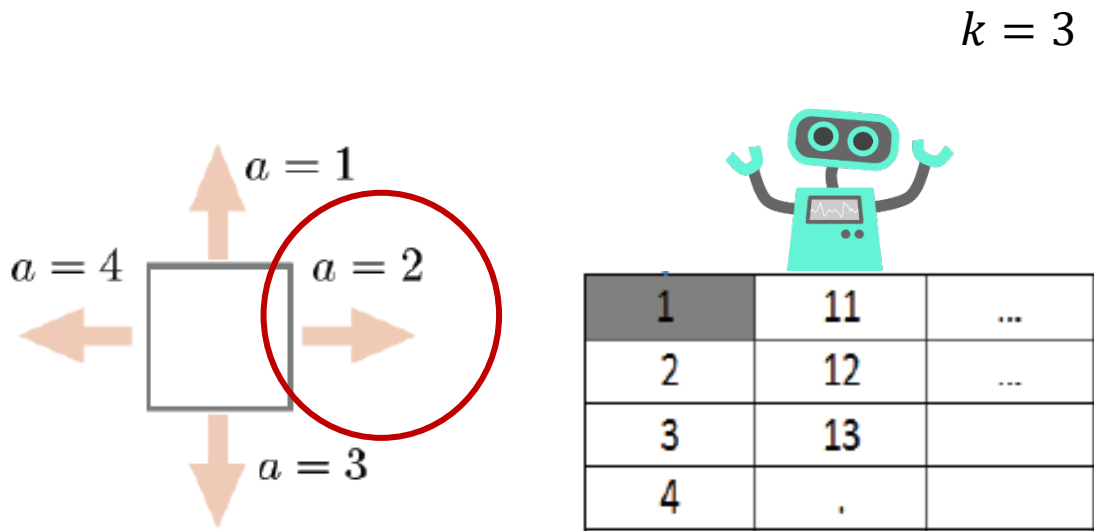


Implementation

Example (continue):

For $k = 3$,

- $s_3 = 11$
- Selected action is $a_3 = 2$
- $s_4 = 21$

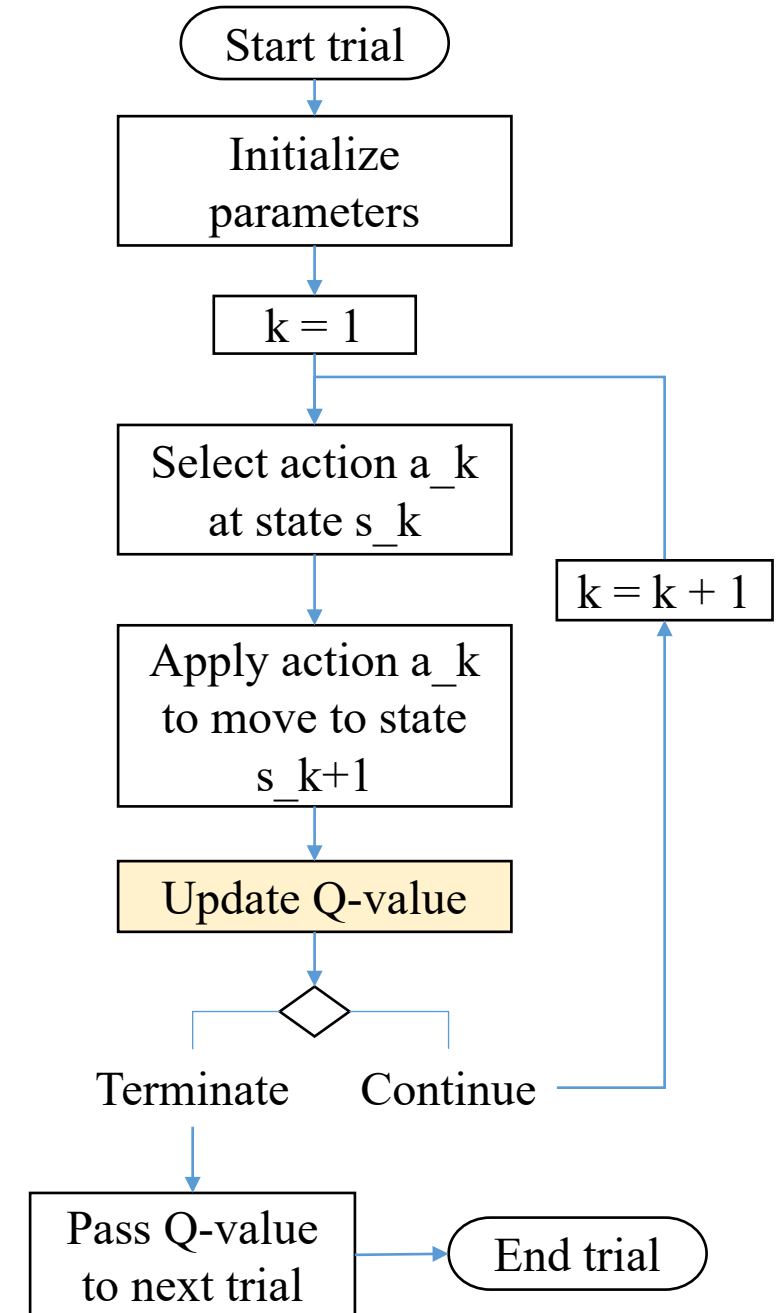
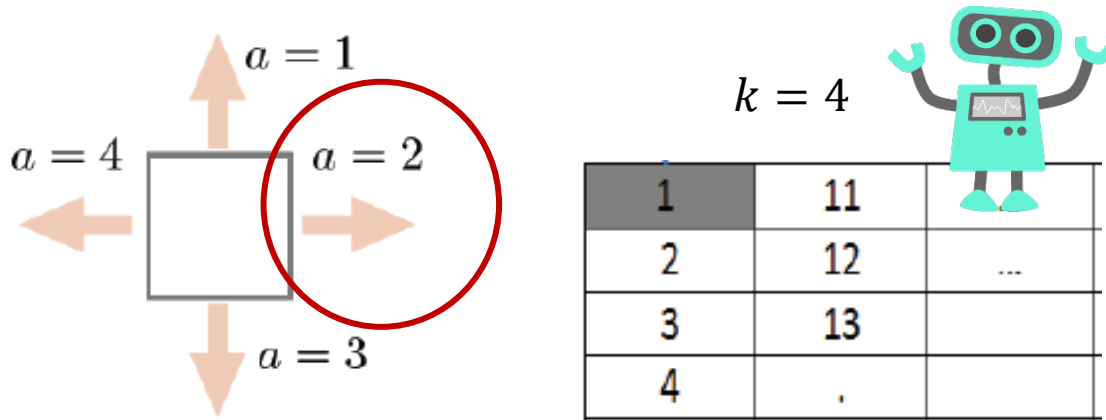


Implementation

Example (continue):

For $k = 3$,

- Receive reward $r_{11,2}$
- $Q_4(11,2) = Q_3(11,2) + \alpha_3(\text{reward}(11,2) + \gamma * \max(Q_3(21,:)) - Q_3(11,2))$



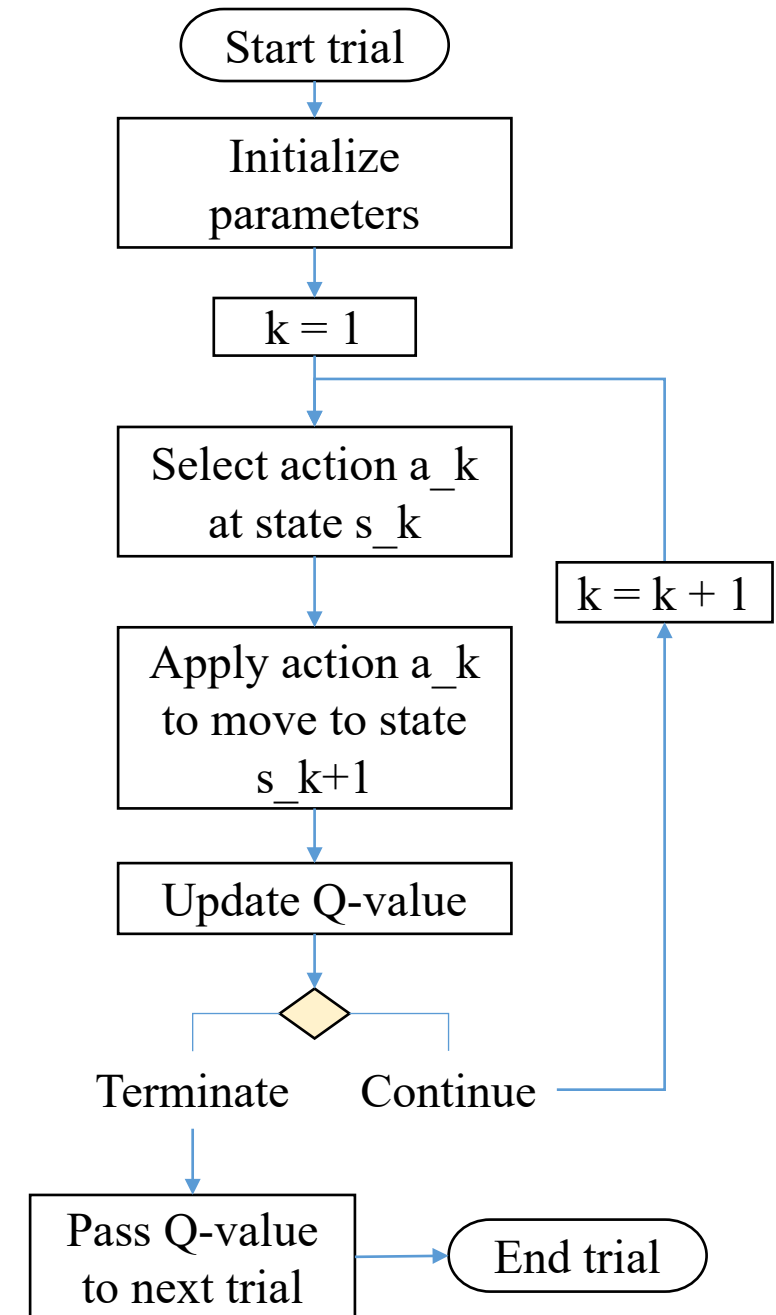
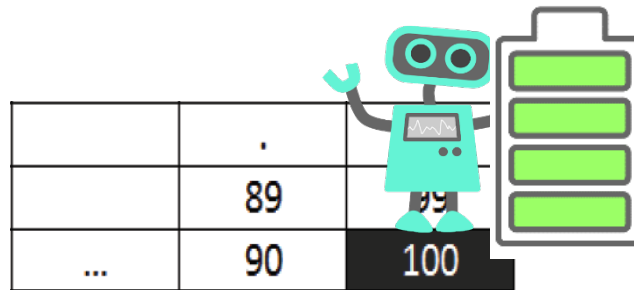
Implementation

Termination condition for each trial:

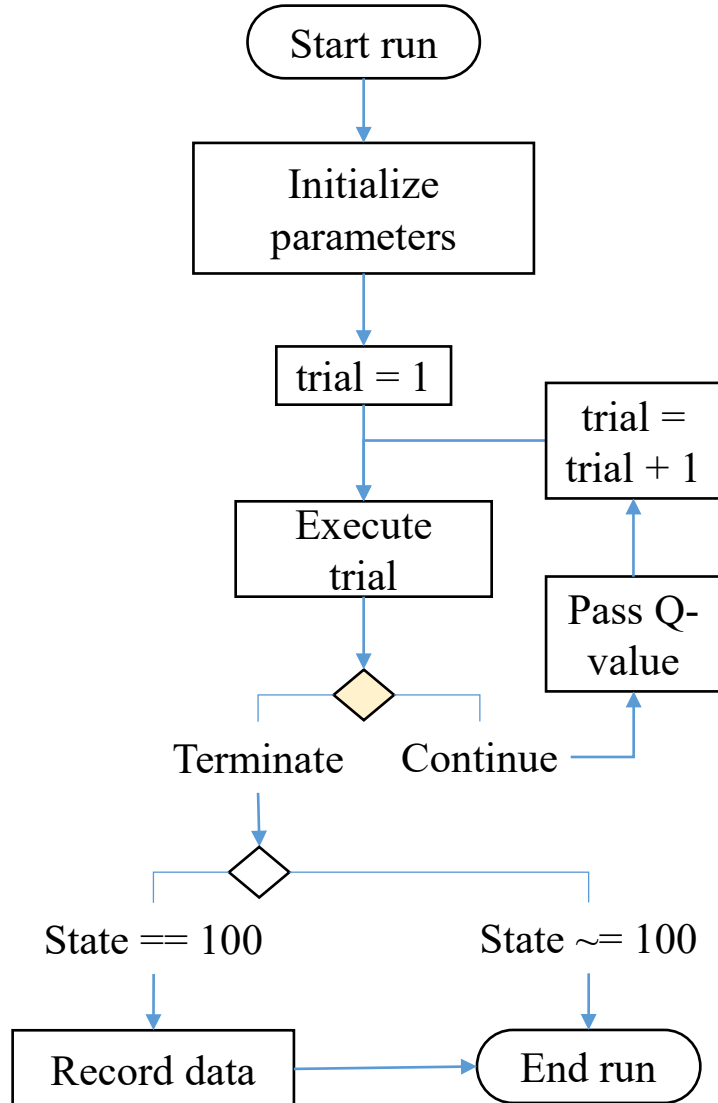
- Robot reaches goal state $s_k = 100$ >> Ideal Case
- $\alpha_k < 0.005$ >> Optional
- Maximum number of time step k_{\max} is reached

Continuation condition for each trial:

- Otherwise



Implementation



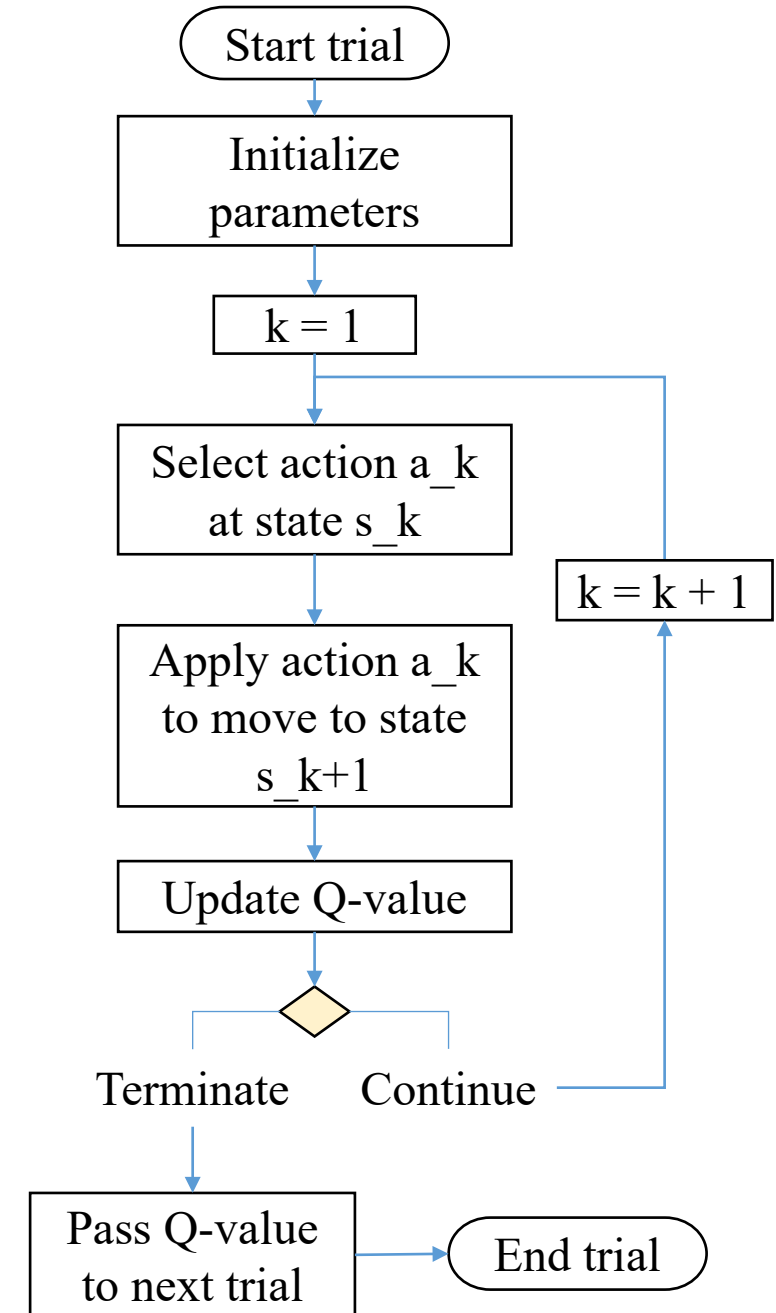
Initial Q-values of next trial is the optimal policy of this trial.

Termination condition for each run:

- Q-function converged to the optimal values >> **Ideal Case**
- Maximum number of trials $trial_{max}$ is reached

Continuation condition for each trial:

- Otherwise





Task I

What to do?

1. Implement Q-learning algorithm in MATLAB
2. Reward function is in task1.mat
- ➔ 3. Discount factor γ and exploration probability ϵ_k are given in Table 1
4. For each set of parameter values,
 - Run 10 runs ($trial_{max} = 3000$ each)
5. Complete report

TABLE I
PARAMETER VALUES AND PERFORMANCE OF Q-LEARNING

| ϵ_k, α_k | No. of goal-reached runs | | Execution time (sec.) | |
|------------------------|--------------------------|----------------|-----------------------|----------------|
| | $\gamma = 0.5$ | $\gamma = 0.9$ | $\gamma = 0.5$ | $\gamma = 0.9$ |
| $\frac{1}{k}$ | ? | ? | ? | ? |
| $\frac{100}{100+k}$ | ? | ? | ? | ? |
| $\frac{1+\log(k)}{k}$ | ? | ? | ? | ? |
| $\frac{1+5\log(k)}{k}$ | ? | ? | ? | ? |

Table 1

For each parameter set

Number of goal reaching runs:

Out of 10 runs, count the runs that the robot ends at state 100.

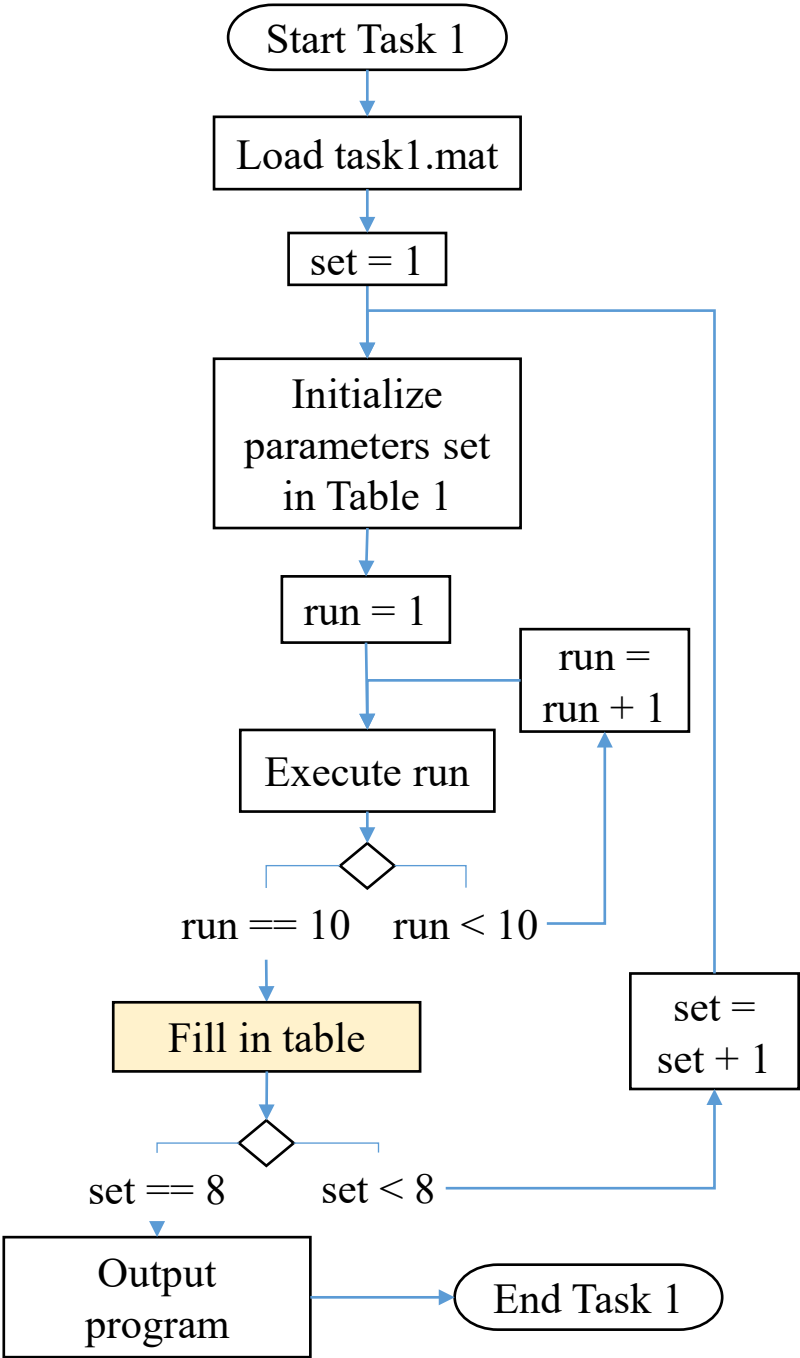
Execution time:

Average the recorded execution time for those goal reaching runs.

TABLE I

PARAMETER VALUES AND PERFORMANCE OF *Q*-LEARNING

| ϵ_k, α_k | No. of goal-reached runs | | Execution time (sec.) | |
|------------------------|--------------------------|----------------|-----------------------|----------------|
| | $\gamma = 0.5$ | $\gamma = 0.9$ | $\gamma = 0.5$ | $\gamma = 0.9$ |
| $\frac{1}{k}$ | ? | ? | ? | ? |
| $\frac{100}{100+k}$ | ? | ? | ? | ? |
| $\frac{1+\log(k)}{k}$ | ? | ? | ? | ? |
| $\frac{1+5\log(k)}{k}$ | ? | ? | ? | ? |

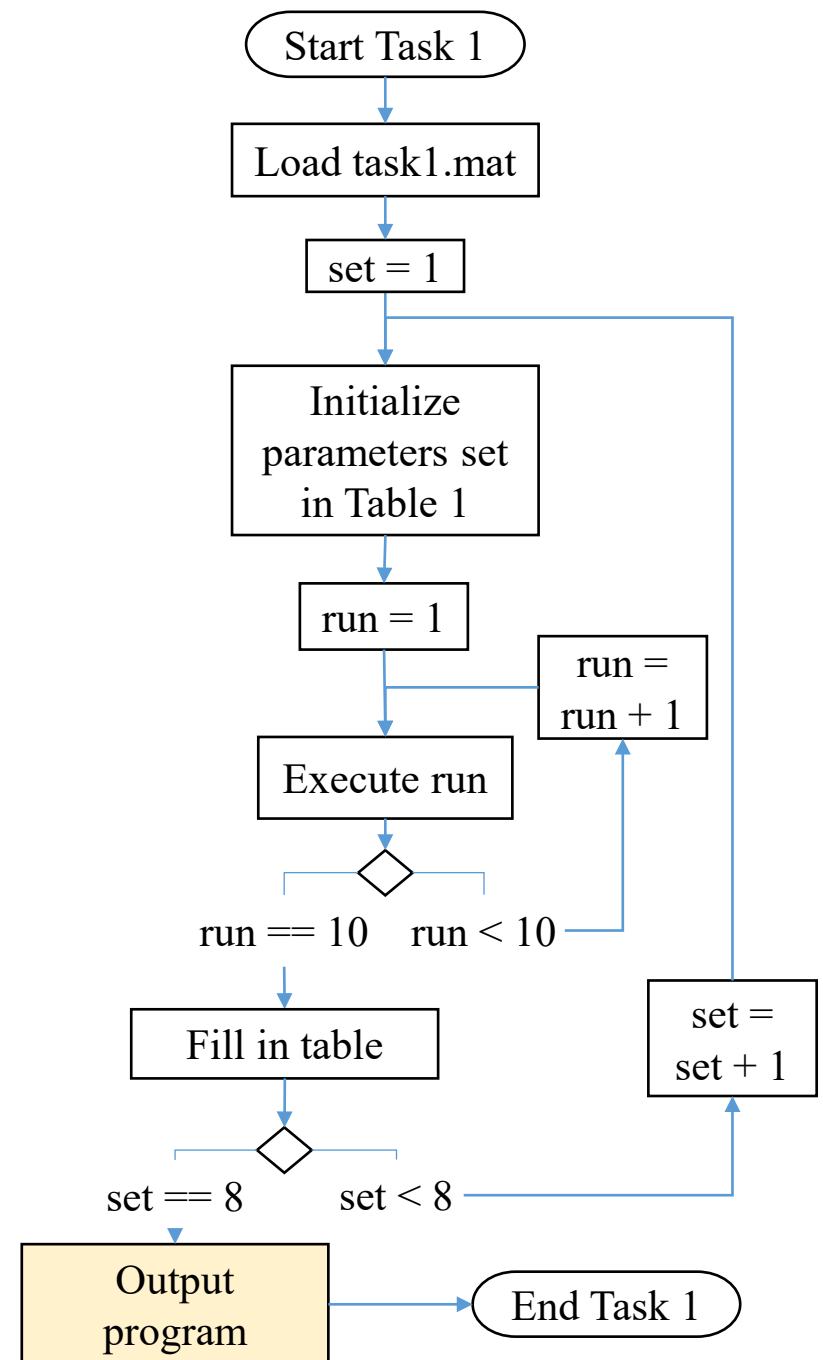


Output

As a column vector, with the position in the column corresponding to a state, and the entry for that position representing the action selected by the optimal policy at that state.

As a 10×10 grid diagram with arrows indicating the action selected by your optimal policy at each state.

As a 10×10 grid diagram showing the (optimal) path taken by the robot as it moves from the initial state to the goal state according to your optimal policy, plus the reward associated with this optimal path.





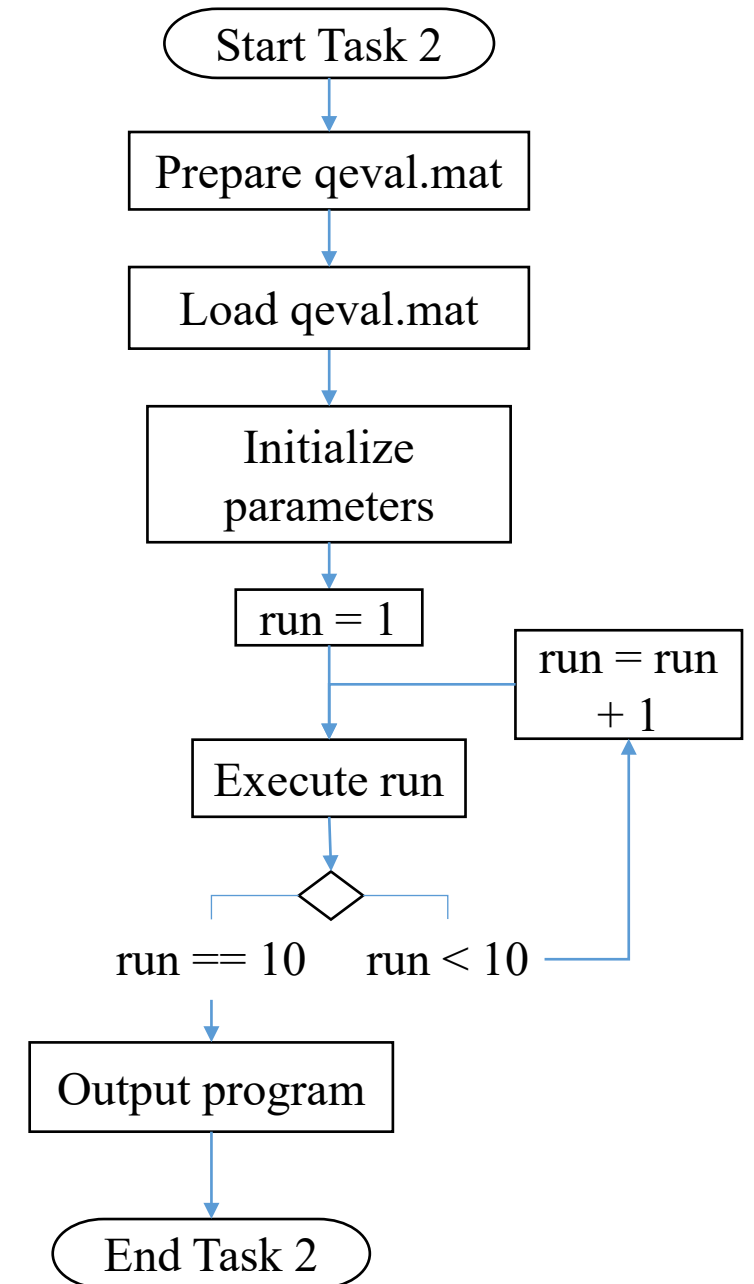
Task II

What to do?

1. Implement Q-learning algorithm in MATLAB.
2. Decide on your discount factor γ and exploration probability ϵ_k .
3. Complete report.

Need to deal with unknown rewards

Note: You can test its execution on your own by making up a qeval.mat file containing dummy sample values.



Assessment

Assessment

The project will be assessed based on the following criteria:

- Comments (with supporting argument) on the results obtained in Task 1
- Presentation. This includes good report style, clarity, and conciseness
- Performance of your M-file program for Task 2 in finding an optimal policy based on the reward function specified in file qeval.mat



Submission

What to submit?

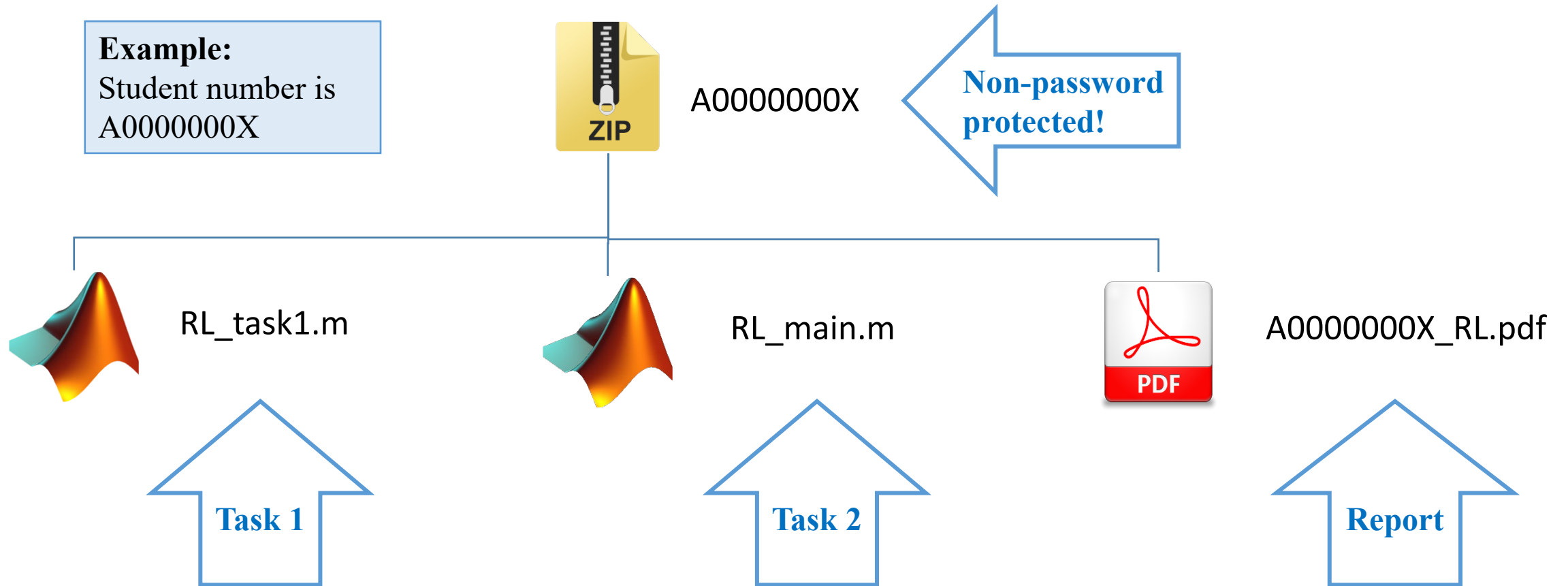
A report (in a **PDF** file) describing the implementation and the results. It must contain a cover page showing:

- Student's name
- Student number
- Student's email address
- Name of the module
- Project title

Name the report as: **StudentNumber_RL.pdf**

What to submit?

And the M-file programs.



Folder: Part II RL project report submission

Report due on 22 April 2022, 23:59 Singapore time

Appendix - MATLAB Course

MATRIX

| | |
|---|-----------------------------|
| Create 2 x 3 matrix | <code>[1 2 3; 4 5 6]</code> |
| Create 4 x 3 matrix of zeros | <code>zeros(4, 3)</code> |
| Find number of rows and columns of matrix A | <code>size(A)</code> |
| Get element at 1 st row and 1 st column of matrix A | <code>A(1, 1)</code> |

REWARD

| | |
|--------------------|---|
| Load 'task1.mat' | <code>load task1.mat</code> |
| Create 'qeval.mat' | <code>save('qeval.mat', 'qevalreward')</code> |

DATA

| | |
|--|--|
| Find the time taken of a block of code | <code>tic</code> <code>% block of code</code> <code>toc</code> |
| Display string with variable | <code>disp(['This is ' variable 'variable.'])</code> |

TRAJECTORY PLOT

Plot coordinate x and y with arrows

- `plot(x, y, '^');` % action 1
- `plot(x, y, '>');` % action 2
- `plot(x, y, 'v');` %, action 3
- `plot(x, y, '<');` % action 4

Set axis min and max

`axis([0 10 0 10])`

Format title

`title(['Execution of optimal policy with
associated reward = ' total_reward])`

Show grid

`grid on`

Start grid from top left corner

`set(gca,'YDir','reverse')`

<https://www.mathworks.com/help/matlab/getting-started-with-matlab.html>

Appendix - Recap

Reward

Total reward for a state transition is given by:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

where,

- R_t determines present value of future rewards

Rewards received k steps in the future is discounted by factor γ^{k-1} .

Small $\gamma \rightarrow$ Focus more on intermediate rewards for next the few steps.

Large $\gamma \rightarrow$ Take into account future rewards more strongly.

Q Function

‘Worth’ of actions at different states is given by:

$$Q^\pi: S \times A \rightarrow \mathcal{R}$$

$$Q^\pi(s, a) = E^\pi[R_t | s_t = s] \rightarrow R_t | s_t = s$$

Deterministic Transition

Expected return from taking action a at state s at time step t by following action π

Optimal Policy

Optimal policy is the state transitions that maximize the Q-values

Slide 164-166

$$Q^{\pi}(s, a) = E^{\pi} [r_{t+1}] + E^{\pi} \left[\gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_t = s \right]$$

Values of Q -function are optimal if they are greater or equal to that of all other policies for all (s, a) pairs, i.e.,

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

Greedy policy

At each s , select a that yields the largest value for the Q -function. When multiple choices are available, such a can be picked randomly

Optimal policy: $\pi^*(s) \in \arg \max_a Q^*(s, a)$

Model-Free Value Iteration

When state transition model is **unknown**, the Q-function can be estimated via iterative update rule by using the reward received from observed state transitions.

$$Q_{k+1}(s_k, a_k) = Q_k(s_k, a_k) + \alpha_k \left(\underbrace{r_{k+1} + \gamma \max_{a'} Q_k(s_{k+1}, a')}_{\text{Estimate of } Q^*(s_k, a_k)} - Q_k(s_k, a_k) \right)$$

Reward of action a at state s

Exploitation:

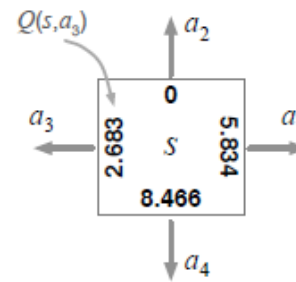
Use greedy policy to select currently known best action

$$a_{k+1} = \max_{a'} Q_k(s_{k+1}, a')$$

Exploration:

Try action other than current known best action

$$a_{k+1} \neq \max_{a'} Q_k(s_{k+1}, a')$$



Exploitation: Take a_4

Exploration: Take a_1, a_2, a_3

ϵ -greedy exploration

Initialize parameters

Input: Discount factor γ ; exploration probability ϵ_k ; learning rate α_k

- Initialize Q -function, e.g., $Q_0 \leftarrow 0$
- Determine the initial state s_0
- For time step k , select action a_k according to:

Select Action

$$a_k = \begin{cases} a \in \arg \max_{\hat{a}} Q_k(s_k, \hat{a}) & \text{Exploitation} \\ & \text{with probability } 1 - \epsilon_k \\ \text{an action uniformly randomly} & \\ \text{selected from all other actions} & \text{Exploration} \\ \text{available at state } s_k & \text{with probability } \epsilon_k \end{cases}$$

Apply Action

Update Q-value

- Apply action a_k , receive reward r_{k+1} , then observe next state s_{k+1}
- Update Q -function with:
$$Q_{k+1}(s_k, a_k) = Q_k(s_k, a_k) + \alpha_k \left(r_{k+1} + \gamma \max_{a'} Q_k(s_{k+1}, a') - Q_k(s_k, a_k) \right)$$
- Set $k = k + 1$ and repeat for-loop for the next time step

Thanks for Listening