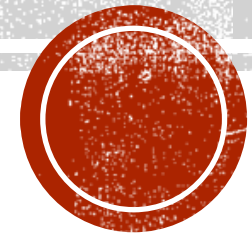# EE5904/ME5404 PART II
# PROJECT 2:
# Q-LEARNING FOR WORLD GRID NAVIGATION
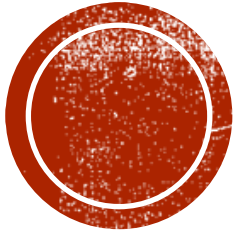
Hanyu Bai

e0816201@u.nus.edu

# PROJECT DESCRIPTION

**Objectives**

1. Competence in implementing the Q-learning algorithm
2. Understanding of the principles of, and implementation issues related to, the Q-learning algorithm.
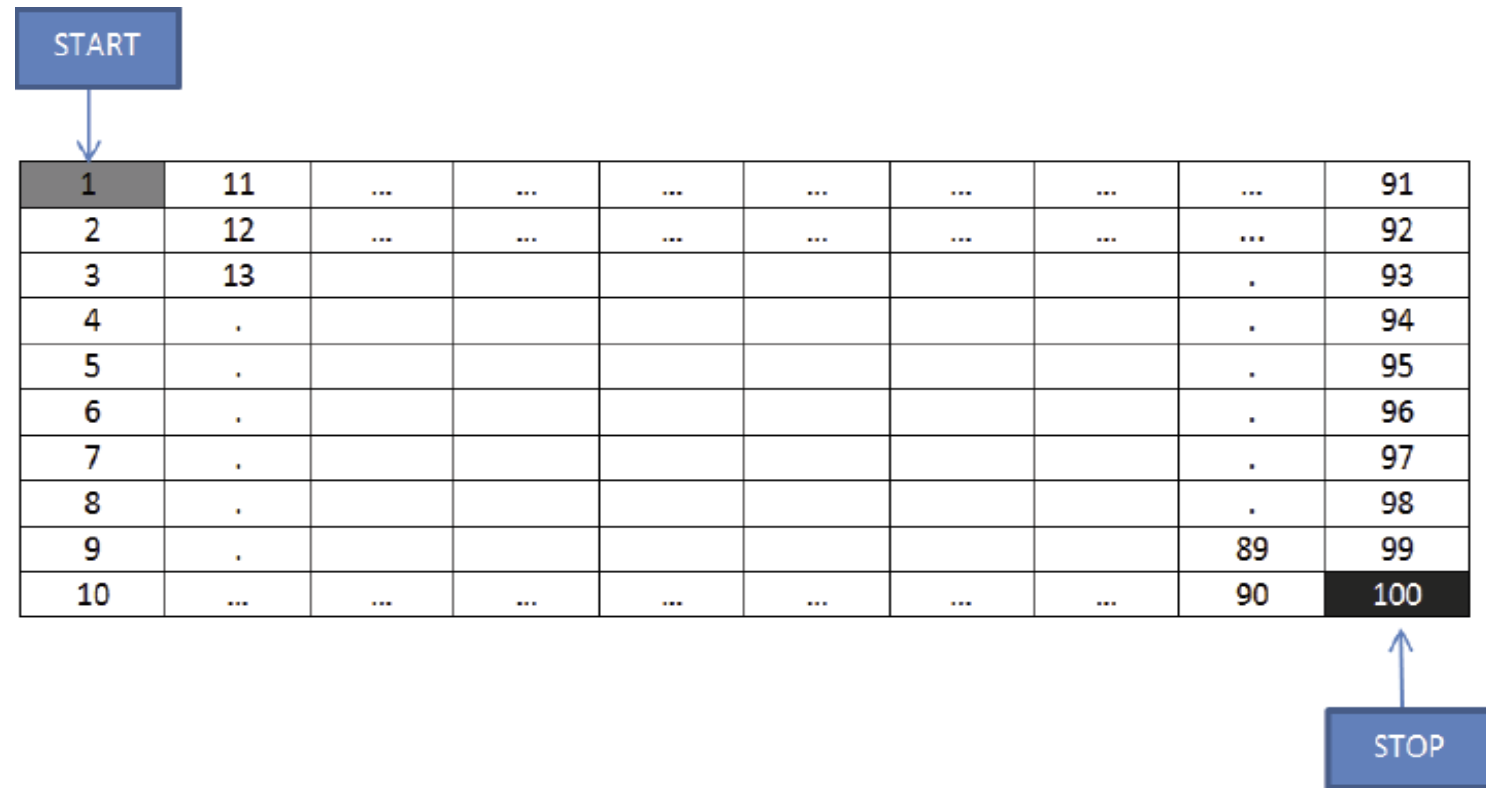
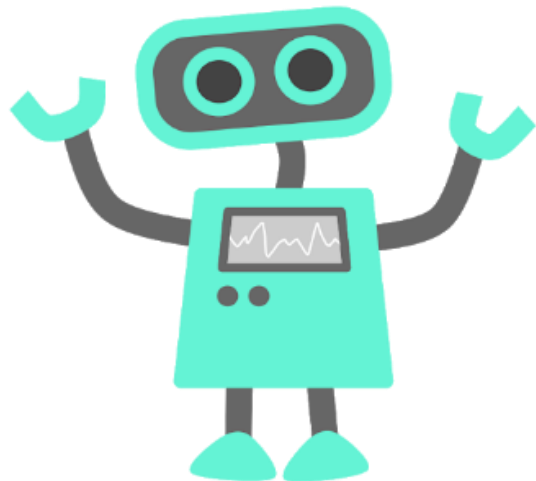**Contents**

1. Tasks
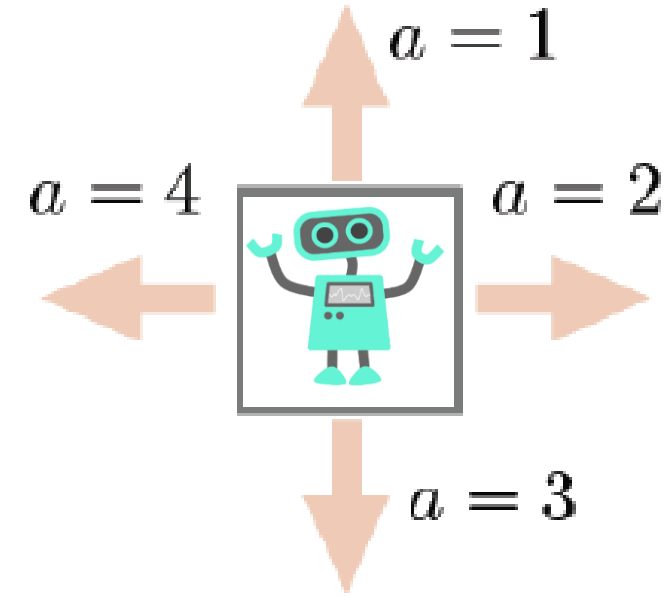2. State transition
3. Reward function
4. Learning

# 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.

# STATE TRANSITION

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

Use dynamic programming methods to find the optimal policy

$$a = 1$$
$$a = 4 \qquad a = 2$$
$$a = 3$$

# REWARD FUNCTION

Files
- Task 1 → "reward" in "task1.mat"
- Task 2 → "qevalreward" in "qeval.mat"

Reward Matrix:
- Column → Action
- Row → State
- 100 x 4

$$
\begin{bmatrix}
r_{1\,1} & \cdots & r_{1\,4} \\
\vdots & \ddots & \vdots \\
r_{100\,1} & \cdots & r_{100\,4}
\end{bmatrix}
$$

Actions

States

# LEARNING

- The robot learns in 1 run, and 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, and 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:**
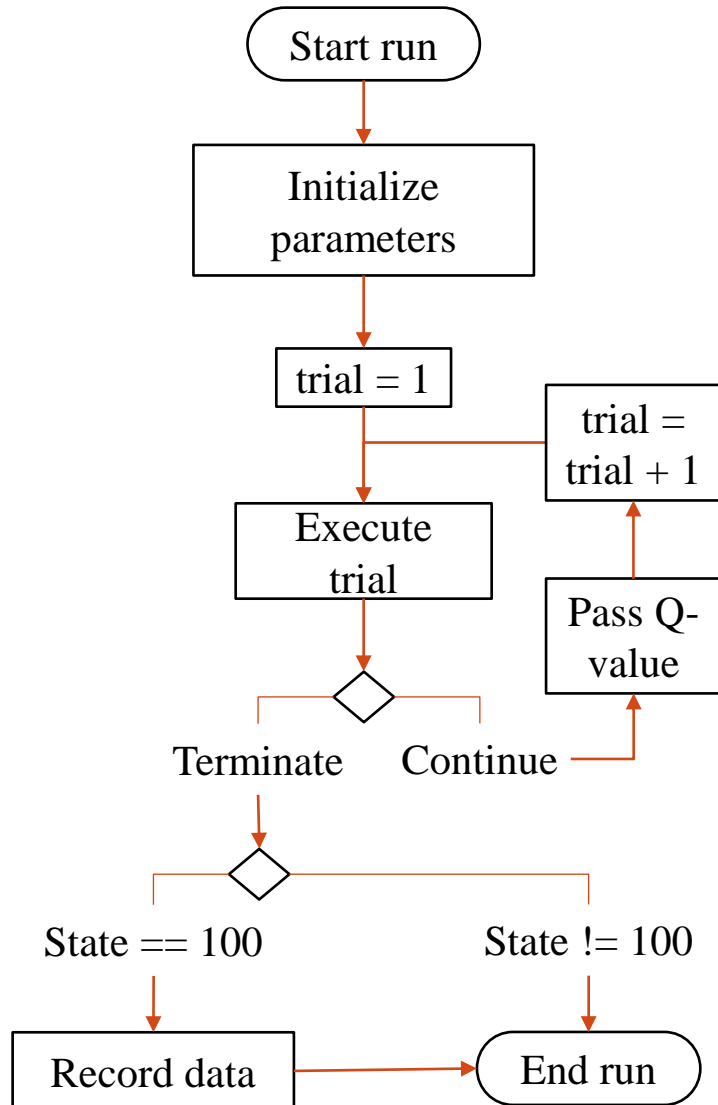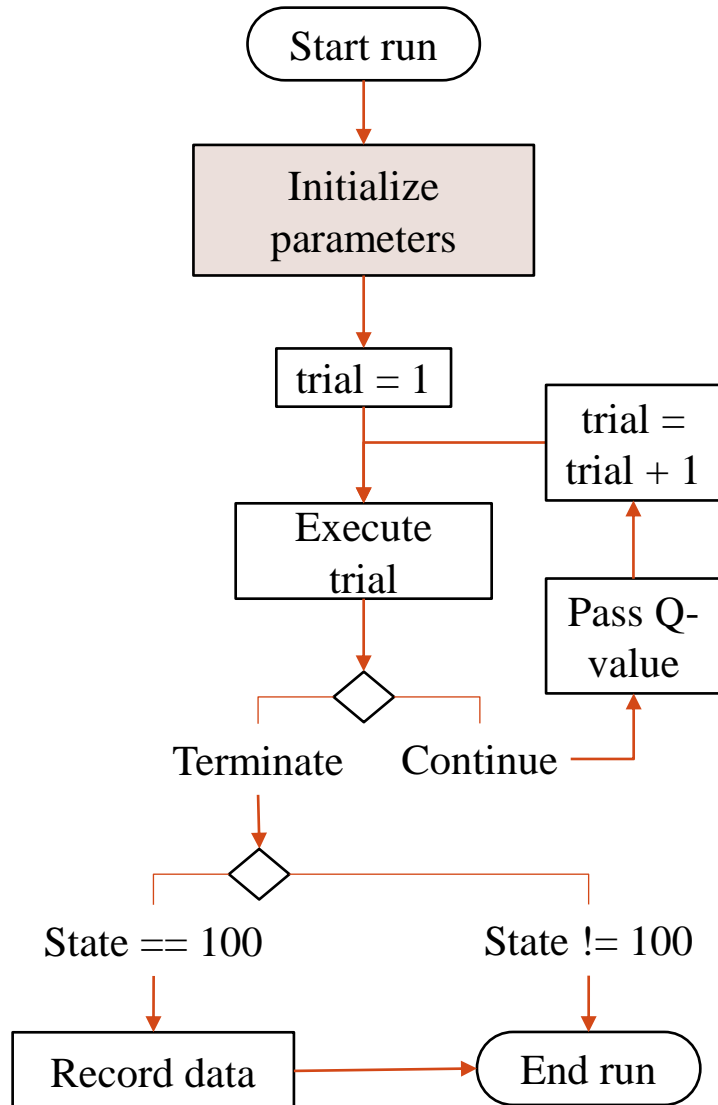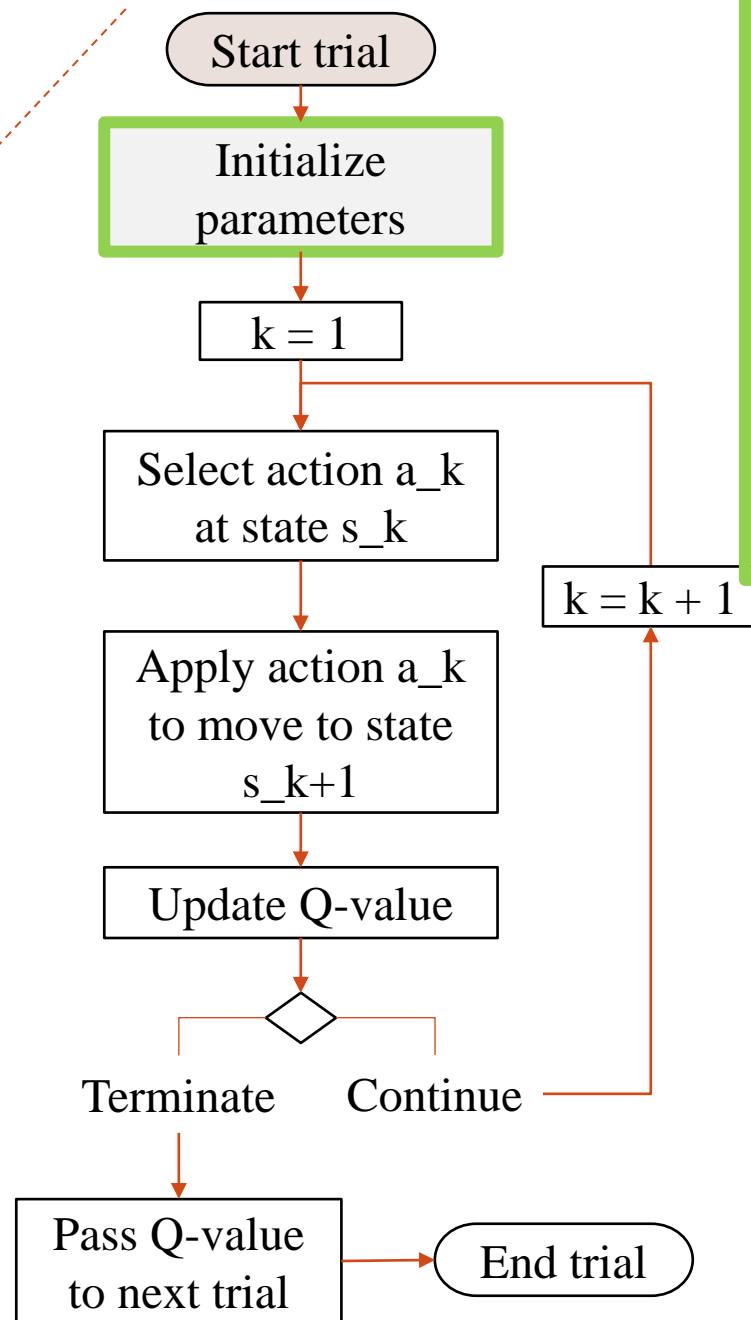
- Initial Q-function $Q_1 \leftarrow 0$ >> Optional
- Trial = 1
- threshold of error of Q-values between trials

# IMPLEMENTATION



Start run → Initialize parameters → trial = 1 → Execute trial → ◇ Terminate / Continue

Continue → Pass Q-value → trial = trial + 1

Terminate → ◇ State == 100 / State ~= 100

State == 100 → Record data → End run

State ~= 100 → End run

Start trial → Initialize parameters → k = 1 → Select action a_k at state s_k → Apply action a_k to move to state s_k+1 → Update Q-value → ◇ Terminate / Continue

Continue → k = k + 1

Terminate → Pass Q-value to next trial → End trial

**Parameters:**
- Discount factor $\gamma$
- Exploration probability $\epsilon_k$
- 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

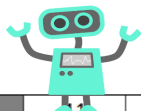**Example:**

For $k = 3$,

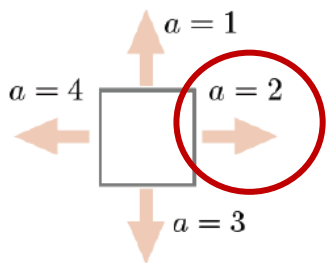- $\epsilon = 1 - \dfrac{1}{k}$

- $s_3 = 11$

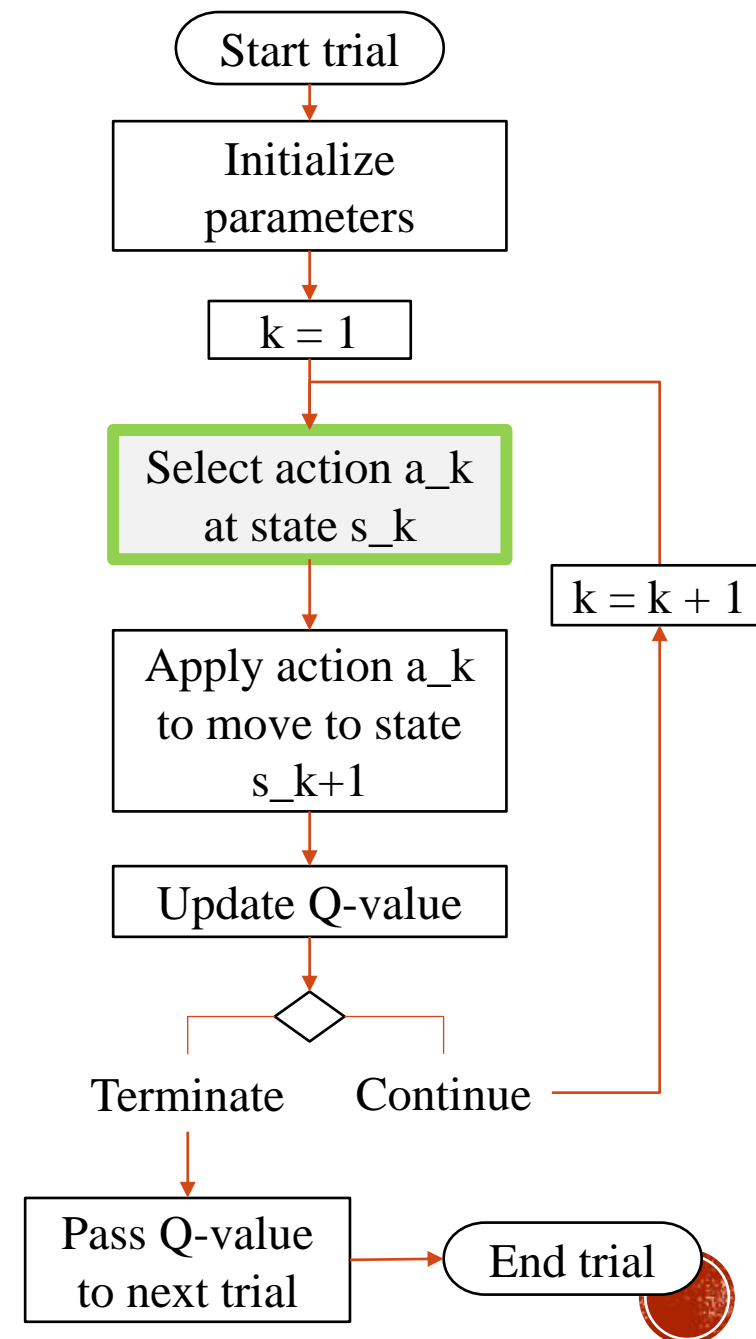| | | | |
|---|---|---|---|
| 1 | | | |
| 2 | 12 | ... | |
| 3 | 13 | | |
| 4 | . | | |

$a_k$

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

Current best action is 2

- Exploitation: $a_3 = 2$ each has $1 - \epsilon = \dfrac{1}{3}$ probability to be selected

- Exploration: $a_3 = 3, 4$ each has $\epsilon = \dfrac{2}{(2)3}$ probability to be selected



$a = 1$

$a = 4$  $a = 2$

$a = 3$

*** $a_3 = 1$ **cannot be** selected due to the **boundary**

Start trial

Initialize parameters

k = 1

Select action a_k at state s_k

k = k + 1

Apply action a_k to move to state s_k+1

Update Q-value

Terminate    Continue

Pass Q-value to next trial

End trial

# IMPLEMENTATION

**Example (continue):**

For $k = 3$,

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



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

# IMPLEMENTATION

**Termination** condition for **<u>each trial</u>**:
- 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 $\gamma$ and exploration probability $\epsilon_k$ 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

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

**For each parameter set in Table 1:**

**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.

Start Task 1

Load task1.mat

set = 1

Initialize parameters set in Table 1

run = 1

run = run + 1

Execute run

run == 10    run < 10

Fill in table

set = set + 1

set == 8    set < 8

Output program

End Task 1

# OUTPUTS

1. 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.

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

3. 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?

***Need to deal with unknown rewards**

1. Implement Q-learning algorithm in MATLAB.

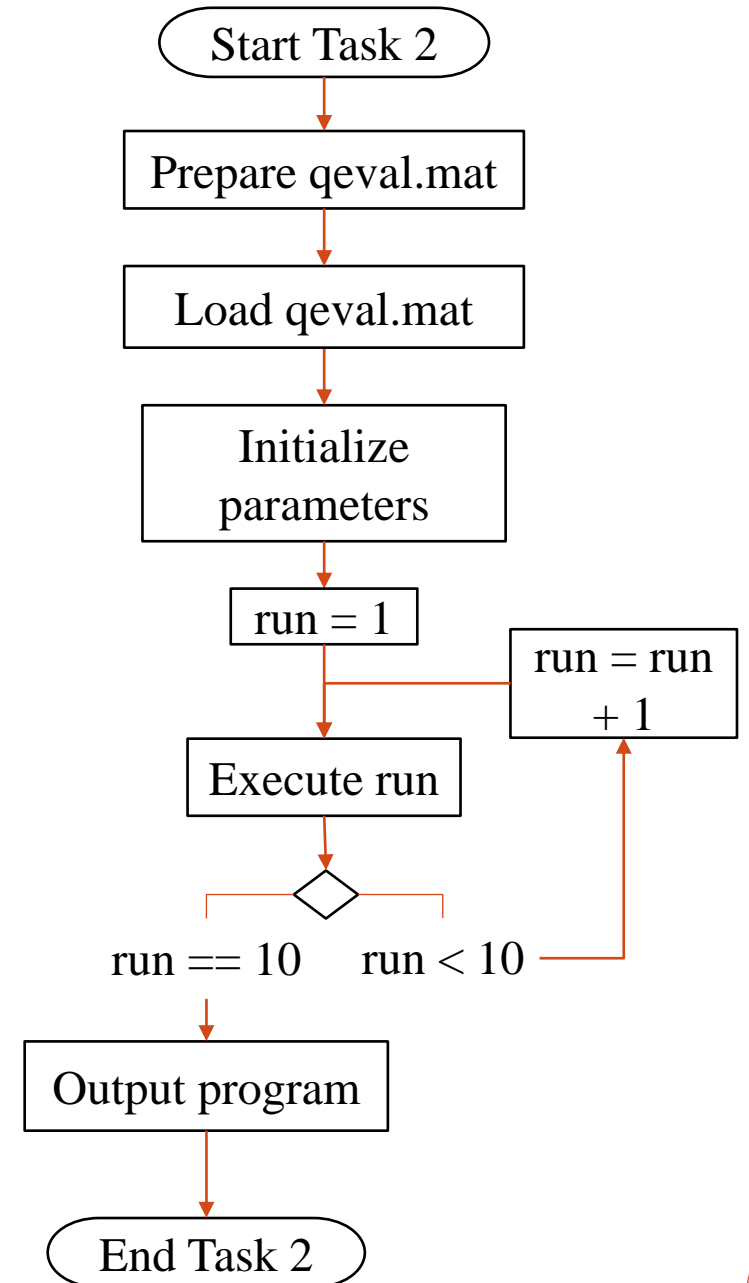2. Decide on your discount factor $\gamma$ and exploration probability $\epsilon_k$.

3. Complete report.

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:**

1. **Comments** (with supporting argument) on the results obtained in Task 1

2. **Presentation of the report**. This includes good report style, clarity, and conciseness

3. **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 with **name** the report as: StudentNumber_RL.pdf) 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

**And** **the M-file programs.**

A0000000X

Non-password protected!

RL_task1.m

Task 1

RL_main.m

Task 2

A0000000X_RL .pdf

Report

# APPENDIX - MATLAB COURSE

## MATRIX

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

## REWARD

| | |
|---|---|
| Load 'task1.mat' | load task1.mat |
| Create 'qeval.mat' | save('qeval.mat', 'qevalreward') |

## DATA

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

# TRAJECTORY PLOT

| | |
|---|---|
| Plot coordinate x and y with arrows | • plot(x, y, '^');   % action 1<br>• plot(x, y, '>');   % action 2<br>• plot(x, y, 'v');   %, action 3<br>• 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') |

# APPENDIX - RECAP

# REWARD

Total reward for a state transition is given by:

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

where,

- $R_t$ determines present value of future rewards

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

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

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

# Q FUNCTION

'Worth' of actions at different states is given by:

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

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

Deterministic Transition

Expected return from taking action $a$ at sate $s$ at time step $t$ by following action $\pi$

# OPTIMAL POLICY

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

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

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.

$$Q_{k+1}(s_k, a_k)$$

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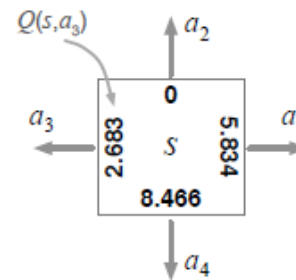
Reward of action $a$ at state $s$

**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

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

$Q(s, a_3)$

$a_2$
$a_3$ 0 $a_1$
2.683 $S$ 5.834
8.466
$a_4$

**Exploitation:** Take $a_4$

**Exploration:** Take $a_1, a_2, a_3$

# $\epsilon$-GREEDY EXPLORATION

**Initialize parameters**

**Select Action**

**Apply Action**

**Update Q-value**

Input: Discount factor $\gamma$; exploration probability $\epsilon_k$; learning rate $\alpha_k$

- 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:

**Exploitation**

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

- 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