

A Deterministic Improved Q-learning for Path-Planning of a Mobile Robot

Amit Konar¹, Indrani Goswami (Chakraborty)¹, Sapam Jitu Singh¹, Lakhmi C. Jain² and Atulya K. Nagar³

¹ETCE Department, Jadavpur University, Kolkata, India.

²University of South Australia, Adelaide.

³Liverpool Hope University, UK.

{konaramit@yahoo.co.in, ic_iit@yahoo.com, sapamjitu@yahoo.com, Lakhmi.Jain@unisa.edu.au, nagara@hope.ac.uk}

Abstract— The paper provides a new deterministic Q-learning with a presumed knowledge about the distance from the current state to both the next state and the goal. This knowledge is efficiently used to update the entries in the Q-table once only by utilizing four derived properties of the Q-learning, instead of repeatedly updating them like the classical Q-learning. Naturally, the proposed algorithm has an insignificantly small time-complexity in comparison to its classical counterpart. Further, the proposed algorithm stores the Q-value for the best possible action at a state, and thus saves significant storage. Experiments undertaken on simulated maze and real platforms confirm that the Q-table obtained by the proposed Q-learning when used for path-planning application of mobile robots outperforms both the classical and extended Q-learning with respect to three metrics: traversal time, number of states traversed, and 90° turns required. Reduction in 90° turnings minimizes the energy consumption, and thus has importance in robotics literature.

Keywords— Agent; Q-learning; Reinforcement learning; Path-planning, Mobile robots.

I. INTRODUCTION

Motion planning is one of the important tasks in intelligent control of a mobile robot. The problem of motion planning is often decomposed into path planning and trajectory planning. In path planning, we need to generate a collision-free path in an environment with obstacles and optimize it with respect to some given criteria [8], [9]. However, the environment may be imprecise, vast, dynamical and partially non-structured [7]. In such environment, path planning depends on the sensory information available at the current state and on the current and cumulative (future) rewards. Since the exact value of the future reward is not known, it is guessed from the knowledge about the robot's world map. The primary advantage of reinforcement learning lies in its inherent power of automatic learning even in presence of small changes in the world map.

Usually, the planning involves an *action policy* to reach a desired goal state, through maximization of a *value function* [1]-[3], which designates sub-objectives and helps choosing the best path. For instance, the value function could be the shortest path, the path with the shortest time, the safest path, or any combination of different sub-objectives. The definition of a task in this context may contain, besides the value function, some *a priori* knowledge about the domains, such as the environmental map, the environmental dynamics and the goal position. The *a priori* knowledge helps the robot generating a plan for motion amidst obstacles, while the lack of such

information of the environment, which might be associated with imprecision and uncertainty. Thus, to have a suitable planning scheme in a cluttered environment, the controller of such kind of robots must have to be adaptive in nature. Several approaches have been proposed to address the problem of motion planning of a mobile robot. If the environment is a known static terrain and it generates a path in advance, it is said to be off-line algorithm. It is called on-line, if it is capable of producing a new path in response to environmental changes.

Machine learning is generally employed in a mobile robot to make it aware about its world map. The early research on mobility management of robots emphasized the needs of supervised learning to train a robot to determine its next position in a given map using the sensory readings obtained by the robot about the environment. Supervised learning is a good choice for mobility management of robots in fixed maps. However, if there is a small change in the robot's world, the acquired knowledge is no longer useful to guide the robot to select its next position. A complete training of the robot with both the old and the new sensory data-action pairs is then required to overcome the said problem.

Reinforcement learning is an alternative learning policy, which rests on the principle reward and punishment. No prior training instances are presumed in reinforcement learning. A learning agent here does an action on the environment, and receives a feedback from the environment based on its action. The feedback provides an immediate reward for the agent. The learning agent here usually adapts its parameter base knowledge obliges the robot to learn it first before invoking the motion planning algorithm.

Our research applies reinforcement learning techniques to real-world robots. Reinforcement learning has been tested in many simulated environments [3] - [11] but on a limited basis in real-world scenarios. A real-world environment poses more challenges than a simulated environment, such as enlarged state spaces [2], increased computational complexity, significant safety issues (a real robot can cause real damage), and longer turnaround times for results. This research measures how well reinforcement-learning technique, such as Q-learning, can be applied to the real robot for navigational problem. In our research, we modified the classical Q-learning algorithm, hereafter called the improved Q-learning for increasing its performance in the path-planning problem.

Performance of a reinforcement learning algorithm is greatly influenced by two important factors used in the control

strategy of the algorithm, popularly known as ‘exploration’ and ‘exploitation’. Exploration usually refers to selecting any action with non-zero probability in every encountered state to learn the environment by the agent. Exploitation, on the other hand, is targeted at employing the current knowledge of the agent to expect achieving good performance by selecting greedy actions [37]. One classical method to balance exploration and exploitation in Q-learning is ϵ -greedy exploration [39], where a parameter ϵ representing exploration probability is introduced to control the ratio between exploration and greedy action selection.

The second alternative method to handle the above problem is to employ Boltzman exploration, where the agent selects an action a with a probability: $\exp(Q(S, a)/T) / \sum_{a'} \exp(Q(S, a')/T)$ in which $Q(S, a)$ denotes the Q-value at state S due to action a , and T is a positive constant called the ‘temperature parameter’. The Boltzman exploration ensures that the agent is more likely to select action a with higher Q-value: $Q(S, a)$. The temperature parameter T adjusts the balance between exploration and exploitation. Large T refers to better exploration, while small T approaching zero indicates almost deterministic (greedy) action selection. Usually these parameters are determined by trial and error to achieve desired performance.

In [38], the authors employed an interesting strategy of internal prediction/estimation to control the balance between exploration and exploitation for efficient adaptability of an agent in a new environment. Here, a reliability parameter RI has been introduced as an internal variable to estimate the ‘expected prediction error’. The variable RI is updated after every action by comparing itself to the actual error. The RI is introduced in the expression of Boltzman exploration by replacing temperature parameters T by $RI=R(S)/\eta$, where η is a positive constant and S denotes the current state of the agent. The factor η is used to normalize the magnitude of $Q(s, a)$ with $R(s)$.

In the context of path-planning by a mobile robot, the robotic agent usually has additional knowledge about the distance from the current state to both the next state and the goal. This knowledge has been efficiently used here during the learning phase of the robotic planner to speed up learning through greedy selection of actions. Four properties (rules) concerning computation of Q-values at state S' from the current estimate of Q-value at state S , where S' is a neighbor of state S , has been developed. Each rule has a conditional part and an action part. If the conditional part is found true, the action part is realized. The conditional part involves checking a locking status of a state along with a distance comparison, where locking of a state indicates that its Q-value needs no further updating. The action part ensures that the Q-value at S' can be evaluated in one step only, and the state S' will also be locked. Thus when conditional part of a property is activated, Q-value at a new state is evaluated once only for ever, and the new state is locked. The proposed improved Q-learning algorithm is terminated when all the states in the workspace are locked.

It is apparent from the above discussion that exploration in the improved Q-learning algorithm takes place when all the rules’ conditional parts do not satisfy at a given situation, whereas exploitation takes place when conditional part of one rule is activated. The balancing of exploration and exploitation is done naturally, and no parameter is involved to control balancing.

The present paper is an extension of the Extended Q-learning (EQL) [17] algorithm, where too the authors presumed that their learning algorithm has knowledge about distance measure of the current state to the next state and the goal. They also listed four properties (rules) of EQL without proof, two of which are derived in the current paper. But the remaining two properties presented here are novel. These two new properties provide alternative conditions for locking, thereby improving relative learning speed of the proposed algorithm in comparison to EQL.

The work presented in this paper is better than the classical [3] and the extended Q-learning by the following counts

1. In classical Q-learning (CQL), Q-values of the states are updated theoretically for infinite number of steps. For all practical purposes however, the algorithm is terminated when the difference in Q-values of each state in two successive iterations is within a prescribed limit. The algorithm is said to have converged under this circumstance, which too requires excessive computational time. The time-complexity of the proposed Q-learning algorithm has been reduced here by locking selected states, where Q-value update is no longer required. The conditions used for identifying the states to be locked are derived here.
2. The extended Q-learning presented in [17] stores only the best action at a state. Naturally, the knowledge acquired by Q-learning in a world map without obstacles cannot be correctly used for planning, particularly when the next state due to the best action is occupied with an obstacle. In the modified Q-learning presented here, the agent is capable of ranking all the actions at a state based on the Q-values at its neighboring states. Consequently, during the planning cycle, if the state corresponding to the best action is occupied with an obstacle, the robot would pickup the next best action. This in one way overcomes one fundamental limitation of the extended Q-learning.
3. Since the extended Q-learning stores only the best action at a state, it cannot take care of the multiplicity of the best actions. Thus if there exist two or more best actions, it selects one of them arbitrarily and saves the selected action with the state. Naturally, the stored action may sometimes involve more turning of the robot during planning than it could have been obtained by an alternative best action. In our algorithm, if we have more than one competitive action at a state during the planning cycle, we select the one ensuring minimum turning of the robot. Thus our present algorithm is energy optimal.

4. Both the extended [17] and the improved Q-learning, algorithms are terminated when all the states are locked. However, the improved Q-learning has 4 locking conditions, including the two of the extended Q-learning. Because of these two additional conditions of locking, the probability of locking of a state in a given interval of time by the improved Q-learning is higher than the extended Q-learning.

The rest of the paper is organized as follows. Classical Q-learning is introduced in section II. Some properties of the Q-learning based on the concept of locking of states are derived in section III. The algorithm for the improved Q-learning is given in section IV. The algorithm for the path planning is given in section V. Computer Simulation is shown in section VI. Experimental details are included in section VII. Conclusions are listed in section VIII.

II. THE CLASSICAL Q-LEARNING

In classical Q-learning, all possible states of an agent and its possible actions in a given state are deterministically known. In other words, for a given agent A, let S_1, S_2, \dots, S_n , be n-possible states, where each state has m possible actions a_1, a_2, \dots, a_m . At a particular state-action pair, the specific reward that the agent acquires is known as *immediate reward*. Let $r(S_i, a_j)$ be the immediate reward that the agent A acquires by executing an action a_j at state S_i . The agent selects its next state from its current states by using a policy. The policy attempts to maximize the cumulative reward that the agent could have in subsequent transition of states from its next state. Let the agent be in state S_i and is expecting to select the next best state. Then the Q-value at state S_i due to action of a_j [36] is given in (1).

$$Q(S_i, a_j) = r(S_i, a_j) + \gamma \text{Max}_d Q(\delta(S_i, a_j), d') \quad (1)$$

where $0 < \gamma < 1$ and $\delta(S_i, a_j)$ denotes the next state due to the selection of action a_j at state S_i . Let the next state selected be S_k . Then $Q(\delta(S_i, a_j), a') = Q(S_k, a')$. Consequently selection of a' that maximizes $Q(S_k, a')$ and in turn $Q(S_i, a_j)$ is an interesting problem.

The classical Q-learning algorithm for deterministic state transitions is given below. The algorithm starts with a randomly selected initial state. An action 'a' from a list of actions a_1, a_2, \dots, a_m is selected, and the agent because of this action receives an *immediate reward* r , and move to the new state following the δ -transition rule given in a table. The Q-value of the previous state due to the action of the agent is

updated following the Q-learning equation (1). Now, the next state is considered as the initial state and the steps of action selection, receiving *immediate reward*, transition to next state and Q-update are represented for ever.

Classical Deterministic Q-learning

For each S, a initialize $Q(S, a) = 0$;

Observe the current state S ;

Repeat

Select $a \in \{a_1, a_2, \dots, a_m\}$ and execute it;

Receive an immediate reward $r(S, a)$;

Observe the new state $S' \leftarrow \delta(S, a)$;

Update the table entry $Q(S, a)$ by

$$Q(S, a) = r(S, a) + \gamma \text{Max}_{a'} Q(\delta(S, a), a');$$

$S \leftarrow S'$;

For ever .

The classical Q-learning requires a memory of (n×m) to keep track of the Q-table. For large n and m, the space complexity thus is high. In the improved Q-learning (IQL) we attempted to reduce the space complexity. The improved Q-learning algorithm presented here involves n Boolean variables called Lock for n states to indicate whether $Q(S, a)$ at state S due to action a needs to be updated. The Lock variables are used to avoid unnecessary update of entries $Q(S, a)$ in the Q-table and thus to save time-complexity.

In IQL, we require n-memories to store n-Lock variables associated with n-states. Here, instead of the Q-table of n×m dimension, we require to store the best Q-value of a state because of any action, and thus require n-memories for n best Q-value of n-states. So, we save some space complexity (nm-2n)=n(m-2).

III. PROPERTIES OF THE IMPROVED Q-LEARNING

This section provides four interesting properties, based on which the IQL algorithm has been developed. The properties stress upon the following two issues.

- 1) If Q-value at any state S (is known and) needs no updating, the state is locked.
- 2) If S is a locked state and S' is a neighboring state of S, and any one of the four properties to be derived is applicable at (S, S') pair, then we can compute Q-value at state S' once only using the property, and the S' is also locked.

If no property is applicable at a given (S, S') pair, state-transition takes place without any updating in Q-value. It may not be out of place to mention here that the first locking in IQL takes place when the agent has a state-transition from the goal state to any of the neighboring states of the goal. The locking of states then is continued as and when one of the four

properties is applicable. The IQL terminates when all the states are locked. We now formally define some parameters to derive the properties.

Let for any state S_k , the distance between the goal state and the next feasible states of S_k are known. Let the next feasible state of S_k be $S \in \{S_a, S_b, S_c, S_d\}$. Let G be the goal and the city block distance between S_a, S_b, S_c, S_d and G be d_{aG}, d_{bG}, d_{cG} and d_{dG} respectively. Let the distance in order be $d_{bG} < d_{aG} < d_{cG} < d_{dG}$. Then the agent should select the next state S_b from its current state S_k . If the Q-value of the state S_b is known, we can evaluate the Q-value of state S_k by the following approach.

$$\begin{aligned} Q(S_k, a') &= r(S_k, a') + \gamma \text{Max}_{a'} Q(\delta(S_k, a'), a'') \\ &= 0 + \gamma \text{Max}_{a'} Q(\delta(S_k, a'), a'') \end{aligned} \quad (2)$$

Now $\delta(S_k, a') = S_a | S_b | S_c | S_d$, where $|$ denotes OR operator.

Therefore,

$$\begin{aligned} &\text{Max}_{a'} Q(\delta(S_k, a'), a'') \\ &= \text{Max}_{a'} Q\{S_a | S_b | S_c | S_d, a''\} \\ &= Q(S_b, a''). \quad (\because d_{bG} < d_{aG} < d_{cG} < d_{dG}) \end{aligned} \quad (3)$$

Combining (2) and (3) we have:

$$\begin{aligned} Q(S_k, a') &= 0 + \gamma Q(S_b, a'') \\ &= \gamma Q(S_b, a''). \end{aligned}$$

Thus if the next state having the shortage distance with the goal is known, and the Q-value of this state is also known, then the Q-value of the current state is simply $\gamma \times$ Q-value of the next state.

Let S_p, S_n and S_G be the present, next and the goal states respectively. Let Q_p and Q_n be the Q-value at the present and the next states S_p and S_n respectively for the best action. Let d_{xy} be the city block distance between the states S_x and S_y . We use a Boolean variable Lock: L_x to indicate that the Q_x value of a state is fixed permanently. We set lock $L_n=1$ if the Q-value of the state n is fixed, and won't change further after L_n is set to 1. The Lock variable for all states except the goal will be initialized to zero in our proposed Q-learning algorithm. We observe four interesting properties as indicated below.

Property 1: If $L_n=1$ and $d_{pG} > d_{nG}$ then $Q_p = \gamma \times Q_n$ and set $L_p = 1$.

Proof: Let the neighborhood state of S_p be $S \in \{S_a, S_b, S_c, S_n\}$, and the agent selects S_n as the next state as $d_{nG} < d_{xG}$ for $x \in \{a, b, c, n\}$.

Now,

$$\begin{aligned} Q_p &= Q(S_p, a) \\ &= r(S_p, a) + \gamma \text{Max}_a Q(\delta(S_p, a), a') \\ &= 0 + \gamma \text{Max}_{a'} Q(S_a | S_b | S_c | S_n, a') \\ &= \gamma \times Q(S_n, a') \quad (\because d_{nG} \leq d_{xG}, \forall x) \\ &= \gamma \times Q_n. \end{aligned} \quad (4)$$

Since $L_n = 1$, and $d_{pG} > d_{nG}$, $\therefore Q_p < Q_n$, and thus $Q_p = \gamma \times Q_n$ for $0 < \gamma < 1$ is the largest possible value of Q_p , and so Q_p should not be updated further. Therefore, $L_p = 1$ is set. \square

Property 2: If $L_p=1$ and $d_{nG} < d_{pG}$ then $Q_n = Q_p / \gamma$ and set $L_n = 1$.

Proof: Since $d_{nG} < d_{pG}$, the agent will select next state n from the current state p . Hence by (4)

$$\begin{aligned} Q_p &= \gamma \times Q_n \\ \Rightarrow Q_n &= Q_p / \gamma. \end{aligned} \quad (5)$$

Since $L_p = 1$, and $d_{nG} < d_{pG}$, $\therefore Q_p < Q_n$, and thus $Q_n = Q_p / \gamma$ for $0 < \gamma < 1$ is the largest possible value of Q_n , and Q_n should not be updated further. So, $L_n = 1$ is set. \square

Property 3: If $L_p = 1$ and $d_{nG} > d_{pG}$ then $Q_n = \gamma \times Q_p$ and set $L_n = 1$.

Proof: If the robot is moving from S_i to the goal G in one step the immediate reward is R , say. On the other hand, if the robot moves from S_i to any state other than the goal G then the immediate reward is zero.

Now, suppose the robot moves from S_p to the goal G in k . Now, as $L_p = 1$, k is the minimum number of state transitions to reach the goal from S_p . So, we obtain

$$\begin{aligned} Q_p &= Q(S_p, a) \\ &= r(S_p, a) + \gamma \text{Max}_{a'} Q(\delta(S_p, a), a') \\ &= 0 + \gamma^k R. \end{aligned} \quad (6)$$

Suppose the robot moves from S_n to the goal G in $k+1$ transition steps. Here, $(k+1)$ is also the minimum number of steps to reach goal G from S_n , failing which the agent would have selected some other state as the next state from the current state S_p .

$$\begin{aligned} Q_n &= Q(S_n, a) \\ &= r(S_n, a) + \gamma \text{Max}_{a'} Q(\delta(S_n, a), a') \\ &= 0 + \gamma^{k+1} R. \end{aligned} \quad (7)$$

Dividing (6) by (7), we have:

$$\begin{aligned}\frac{Q_p}{Q_n} &= \frac{\gamma^k R}{\gamma^{k+1} R} \\ &= \frac{1}{\gamma}\end{aligned}$$

$$\Rightarrow Q_n = \gamma \times Q_p. \quad (8)$$

Since $d_{nG} > d_{pG}$, $Q_n < Q_p$. So, $Q_n = \gamma \times Q_p$ has largest possible value for $0 < \gamma < 1$. Further, as $L_p=1$, and S_n is the nearest state to S_p with respect to given distance metric, therefore L_n is set to 1. \square

Property 4: If $L_n=1$ and $d_{nG} > d_{pG}$ then $Q_p = Q_n/\gamma$ and set $L_p = 1$.

Proof: Since $d_{nG} > d_{pG}$, Q_n can be evaluated from Q_p by (8). Thus

$$\begin{aligned}Q_n &= \gamma \times Q_p \\ \Rightarrow Q_p &= Q_n / \gamma.\end{aligned} \quad (9)$$

Since $L_n = 1$, and $d_{nG} > d_{pG}$, $\therefore Q_n < Q_p$, and thus $Q_p = Q_n / \gamma$ for $0 < \gamma < 1$ is the largest possible value of Q_p , and Q_p should not be updated further. So, $L_p = 1$ is set. \square

IV. THE IMPROVED Q-LEARNING ALGORITHM

The classical Q-learning employs a Q-table to store the $Q(S, a)$ for $S=S_i$ to S_n and $a=a_j$ to a_m . Thus it requires an array of $(n \times m)$ size. In the improved Q-learning, we, however, require to store only the Q-values at a state S for the best action. Thus for n states, we need to store n Q-values. Besides the Q-storage, we in addition require n Boolean Lock variables, denoted by L_i for state S_i , $i=1$ to n, depicting the current status of the state. If the Lock-variable at a state is 1, then the Q-value at that state need not be updated further.

In path-planning application of mobile robots, the environment can be partitioned into non-overlapped grids, called states. Thus a state can have four neighbors. Consequently, during the planning phase, the robot can determine the best action to move to the next (best) state S_n from the current S_p by identifying the neighboring state having the largest Q-value. The corresponding action of the robot to move to the next (best) state S_n is apparent.

The proposed algorithm for improved Q-learning has 2 main steps; i) initialization, ii) Q-table updating. In the initialization phase, the lock variable at all states except the goal state S_G is set to zero. The immediate reward from any neighboring state to the goal state is set to 100. The discounting factor γ and the initial state are fixed up.

In the present update policy of the Q-table, if $L_p(L_n)$ is 1, then $L_n(L_p)$ will be set to 1. However, in the initialization phase only L_G is set to 1. Thus, unless the current or the next state = L_G , there will be no update in the Q-table. In order to have L_p or $L_n = L_G$, the robot usually has to wander in its world map for a finitely large number of iterations. To avoid the

unnecessary execution of the Q-table update, we add a small repeat-until loop between the two main phases of the program. This loop continues selecting an action and execute it (without updating Q-table) until the robot reaches the goal.

Once the robot reaches the goal, the first repeat-until loop exits, and the Q-table updating is initiated. The process of Q-table updating is continued until all the states are locked.

The pseudo code of the improved Q-learning is given below.

Pseudo Code for Improved Q-learning

1. Initialization

For all $S_i, i=1$ to n except $S_i = S_G$

{set $L_i = 0; Q_i = 0;$ }

L_G (for goal S_G) = 1;

Q_G (for goal S_G) = 100;

Assign γ in (0, 1) and initial state = S_p ;

2. Repeat

{

Select a_i from $A = \{a_1, a_2, \dots, a_m\}$ and execute it;

} **Until** $S_p = S_G$;

3. Update Q - table:

Repeat

{

a) Select a_i from $A = \{a_1, a_2, \dots, a_m\}$;

b) Determine d_{nG} and d_{pG} ;

If ($d_{nG} < d_{pG}$)

Then if ($L_n = 1$)

Then if ($L_p = 0$)

Then $\{Q_p = \gamma \times Q_n; L_p = 1;\}$

Else if ($L_p = 1$)

Then $\{Q_n = Q_p / \gamma; L_n = 1;\}$

Else if ($L_p = 1$)

Then if ($L_n = 0$)

Then $\{Q_n = \gamma \times Q_p; L_n = 1;\}$

Else if ($L_n = 1$)

Then $\{Q_p = Q_n / \gamma; L_p = 1;\}$

} **Until** $L_i = 1$ for all i without obstacle;

Theorem 1: The entries in the Q-table for the best action in the classical Q-learning have the same value as in the improved Q-learning.

Proof: Let Q_p and Q_n be the Q-value for the best action at the present state S_p and the next state S_n . In classical Q-learning Q_p is updated, when the agent has a transition from S_p to S_n . However, Q_p would attain maximum value in a learning epoch, if the Q_n had already attained the maximum value.

Now, if S_p is closest to the goal, then $Q_p = \gamma \cdot R$, where R is immediate reward and γ is the discounting factor. So, if S_p is at a distance of k through a shortest path measured by city block distance, then $Q_p = \gamma^k \cdot R$. Q_p if updated later cannot exceed $\gamma^k \cdot R$, as it is at a shortest distance k w.r.t. the goal.

In the improved Q-learning, if $d_{pG} > d_{nG}$, then by property (1) and (2) $Q_p = \gamma \cdot Q_n$. Now, if $L_n = 1$, i.e., state n is at a

shortest distance (k-1) to the goal, $Q_n = \gamma^{k-1} \cdot R$, and then $Q_p = \gamma \cdot Q_n = \gamma \cdot \gamma^{k-1} \cdot R = \gamma^k \cdot R$. So, when $d_{pG} > d_{nG}$ and $L_n = 1$, $Q_p = \gamma^k \cdot R$, and this value of Q_p should not change further.

Further, if $d_{nG} > d_{pG}$, then by property (3) and (4) $Q_p = Q_n / \gamma$. Now, if Q_p is at a shortest distance (k+1) from the goal, then $Q_n = \gamma^{k+1} \cdot R$ and $\therefore Q_p = Q_n / \gamma = \gamma^k \cdot R$, and $\gamma^k \cdot R$ is the largest possible value of Q_p .

In the improved Q-learning, irrespective of $d_{pG} > d_{nG}$ or $d_{nG} > d_{pG}$ it is found that $Q_p = \gamma^k \cdot R$. So, both the classical Q-learning and the improved Q-learning have a steady state Q-table with $Q_p = \gamma^k \cdot R, \forall p$. \square

We now determine the space- and time-complexity of the Improved Q-learning.

Space-Complexity: In classical Q-learning, if there are n states and m action per state, then the Q-table will be of (m×n) dimension. In the Improved Q-learning, 2 storages are required for each state, one for storing Q-value and other for storing value of the lock variable of a particular state. Thus for n number of states, we require a Q-table of (2×n) dimension. The saving in memory in the present context with respect to classical Q thus is given by $mn - 2n = n(m - 2)$, which is of the order of mn.

Time-Complexity: In classical Q-learning, the updating of Q-values in a given state requires determining the largest Q-value, in the next state for all possible actions by Equation (1). Thus if there are m possible actions at a given state, maximization of m possible Q-values, require m-1 comparison. Consequently, if we have n number of states, the updating of Q values of the entire Q table by classical method requires $n(m - 1)$ comparisons. Unlike the classical case, here we do not require any such comparison to evaluate the Q values at a state S_p from the next state S_n . But we need to know whether state n is locked that is, Q-value of S_n is permanent and stable. Thus if we have n number of states, we require n number of comparison. Consequently, we save a time $n(m - 1) - n = nm - 2n = n(m - 2)$, which is of the order of mn.

V. THE PATH PLANNING ALGORITHM

The Q-learning algorithm presented above stores the Q-values at each state for the best action. After the learning is completed, i.e., all the states are locked, the Q-table can be used for path planning application. During path planning, the robot while at state S_p identifies the next best state S_n , where the Q-value is higher than the Q-value of other neighboring states of S_p . However, if there exist more than one next state having the largest Q-value among the neighboring state of S_p , the robot ideally would select any one of them.

In this paper, we, however, economically select the next state S_n , while the robot is at S_p , based on the torque requirement. For example, let there exist two states S_{n1} and S_{n2} having the largest Q-value around the neighbor of S_p . Then the robot would select S_{n1} as the next state, if the angular rotation to move to S_{n1} , is smaller than that of S_{n2} . A small angular turning requires less torque to be generated by the robot. Consequently, the robot in the proposed planning algorithm consumes minimum energy as the torques generated during successive movement of the robot toward the goal is optimally selected.

Let the present state be S_p and S_n be the next state with the largest Q-value. S_r be a next state of S_p , and S_G be the goal state. We now develop a path-planning algorithm to determine an obstacle-free trajectory of optimal path-length and energy for the robot between an arbitrary starting point and a fixed goal point in the pre-trained world map.

Pseudo Code for Path Planning

BEGIN

1. $CURRENT \leftarrow S_p$;

2. If $Q_n > Q_r, \forall r$

Then $\forall n$ select n' from n , such that angular rotation required to reach n' is minimum and n' is obstacle - free;

3. Go to $NEXT$;

4. $CURRENT = NEXT$;

5. Repeat from step 2 until $NEXT = GOAL$;

END

VI. COMPUTER SIMULATION

In our computer simulation, we consider an environment of 20×20 grids, where each grid is given a state number. For example, a grid located at position (x, y) defined in the Cartesian coordinate reference frame has a

$$stateno = (x - 1) \times rowsize + y \quad (10)$$

where $rowsize$ denotes the number of grids in a row.

Performance of the Q-learning algorithm has been studied here in two phases. First, a given world map is trained by the proposed Q-learning algorithm. Second, the trained world map with a known Q-table is used to generate a trajectory of motion of the robot between an arbitrarily selected initial position and the fixed goal position in the said world map. The performance metrics used here to compare the relative performance of our proposed algorithm with the classical Q- and the extended Q-learning [5] include the convergence time of the learning algorithm, total time taken to execute a plan of motion in a pre-trained world map, and the number of $\pm 90^\circ$ rotations involved to completely execute the plan. While convergence is considered for the learning algorithm, the issues of time and energy consumptions, the latter being measured in terms of $\pm 90^\circ$ turnings, are part of the planning algorithm. The performance analysis considers both the learning and the planning algorithms together as both of them jointly determine the overall performance in the path-planning application.

The performance of the Q-learning algorithm is studied under four experimental settings. First, the training and planning is performed on the same obstacle-free world maps. Second, the training is performed on obstacle-free map, but before planning algorithm is executed obstacles are added in the map. Third, both training and planning are performed on the same map with few obstacles. Lastly, the training was performed in a map with few obstacles, and a few more obstacles are added before planning. The classical Q-, the extended Q- and the improved Q-learning algorithms are compared with the well-known Dijkstra's shortest path finding algorithm in a graph and the heuristic A* algorithm under the above four experimental settings.

A. Experiments

Let S and G be the starting and the goal states in all the world maps considered for path-planning in the following experiments. In each experiment, the Q-table is obtained by three different Q-learning algorithms: the classical Q-, the extended Q- and the improved Q-learning. After the learning is over, the robot is kept at the starting position with the heading direction in the east. The experiments and the corresponding results are briefly outlined below.

Experiment 1: The first experiment is carried out on a world map of 20×20 grids as shown in Fig. 1 to compare the relative performance of the three distinct Q-learning algorithms. The classical Q-learning with random action selection and Boltzman action selection with $T=0.01$ have found to converge after 50026 and 8276 iterations.

The extended and the improved Q-learning are executed until all the Lock variables associated with each state are set to 1. The numbers of iterations required to learn the world map of Fig. 1 by the extended Q-learning was found to be 20273, whereas the improved Q-learning requires 6502 iterations to set all the states locked.

The experiment 1 reveals that the paths obtained by the planning algorithm by acquiring knowledge in the Q-table by all the three algorithms are optimal (having 21 state transitions) with respect to a measure of city block distance. However, the improved Q-learning requires minimum turning and thus consumes minimum energy to execute the complete task of path-planning (Table- I). As A* and Dijkstra's algorithms search optimal paths in real time, they respectively consume almost double and 4-times the time required for path-planning by IQL.

Experiment 2: The second experiment is concerned with training in the obstacle-free world map of Fig. 1 by the three Q-learning algorithms. After the training is over, we add obstacles in the map and change the starting position as indicated in Fig.2 to Fig.5, and execute the respective planning programs. The resulting paths obtained by the planning algorithms with the acquired knowledge stored in the Q-tables by respective Q-learning algorithm reveal interesting observations. First, the resulting Q-tables obtained by the classical and the improved Q-learning based techniques help the planning algorithms to construct optimal paths in all the

maps of Fig. 2-5. Fig. 2, 4 and 5 exhibit an immature termination of the trajectory of motion by the robot, when the Q-table is updated by the extended Q-learning. This happens because the extended Q-learning algorithm stores only the best action at each state, which sometimes is occupied by an obstacle.

The scenario, however, is different in case of modified Q-learning. In modified Q-learning, the agent during execution of the plan determines the best neighborhood state having the largest Q-value. If the neighboring state thus selected is occupied by an obstacle, the robot selects the next feasible neighboring state with the largest Q-value, and selects it as the next state. Consequently, even only one neighboring state, which was the previous state of the agent is unoccupied with obstacle, the agent can return to that state. Thus the planning algorithm never gets stuck to a state as in case of the extended Q-learning.

It is noteworthy that 90° turns taken by A* and Dijkstra's algorithms are smaller in comparison to that by IQL in Fig. 4 and 5. The justification of the results is due to the phenomenon that the Q-values stored do not carry information about turning angles. So a search algorithm looking for a shortest path naturally identifies a trajectory with less turning angles.

Experiment 3: The third experiment is carried out in world map (Fig. 6) with obstacles during both the learning and the planning phase. The number of learning epochs required for convergence by different algorithms are 50604 for classical Q-learning with random action selection, 9104 for classical Q-learning with Boltzman action selection ($T=0.01$), 21701 for EQL and 6546 for IQL.

The planned trajectories obtained by consulting respective Q-tables for CQL, EQL and IQL algorithms are shown in Fig. 6. The paths planned by Dijkstra's shortest path finding and A* algorithms are also included in Fig. 6 for comparison. It is apparent from the figure that the shortest path (no. of state transitions=13) is obtained, when the Q-table is updated both by the extended Q-learning and the improved Q-learning. On the other hand, the path constructed by consulting the Q-table obtained by the classical Q-learning is excessively longer with 19 state transitions. The torque requirement is the smallest in the improved Q-learning, as the number of turnings required here is minimum (once only). A* and Dijkstra's algorithms take excessively large planning time.

Experiment 4: The last experiment is concerned with training in a map with 5 dark obstacles, and planning trajectory in that map with 3 additional shaded obstacles (Fig. 7). It is apparent from Fig. 7 that the improved and the classical Q-learning induced Q-table help the robot generate complete trajectories of motion of the robot. However, the planning algorithm realized with the Q-table obtained by the extended Q-learning fails to generate a complete trajectory. The number of state transitions and torque requirement are also minimum in case of improved Q-learning. In Fig. 8, we consider 5 dark obstacles during the training and 4 additional shaded obstacles during the planning phase. It is noted from the figure that minimum number of state transitions take place in case of path planning using the Q-table obtained by improved Q-learning.

The turning required by the said trajectory is also minimum when compared to the other trajectories.

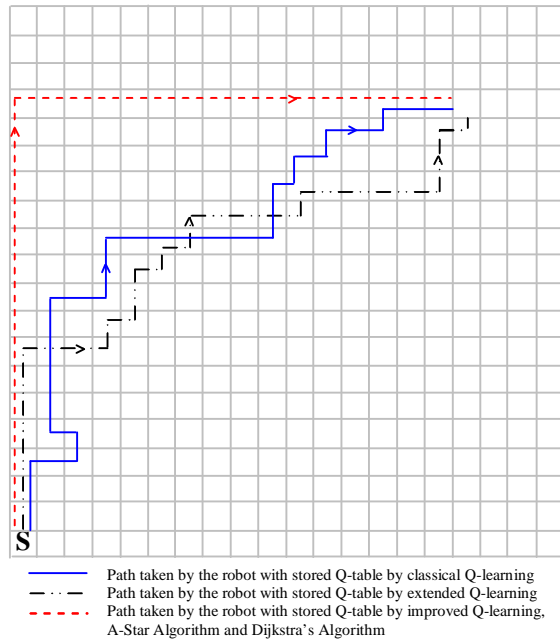


Fig. 1: World map 1 without obstacle. Paths taken by the robot with stored Q-table by the different algorithms are shown in the fig.

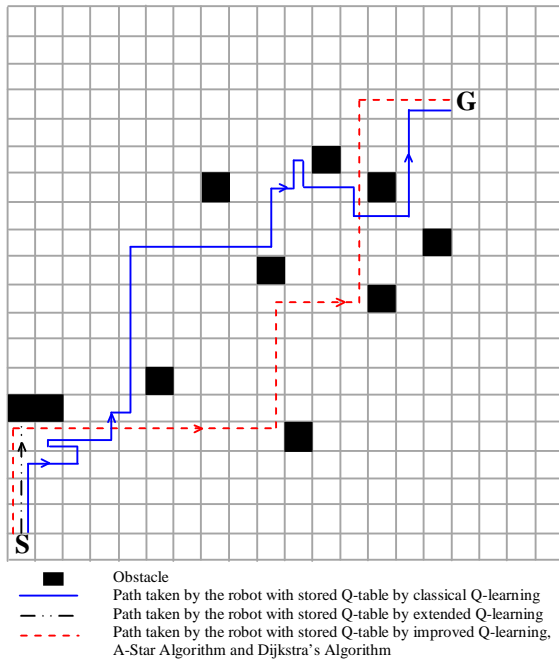


Fig 2: World map 2 with obstacles. Paths taken by the robot with the stored Q-table by the different algorithms are shown in the fig. The robot with the stored Q-table by the extended Q-learning fails to reach the goal.

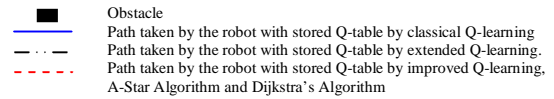
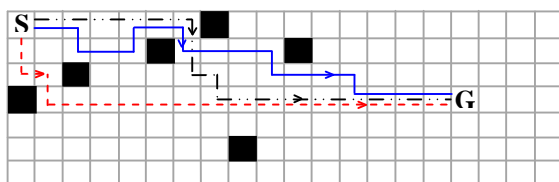


Fig 3: A Partial World map 3 with obstacles. It shows the paths taken by the robot with the stored Q-table of different algorithms.

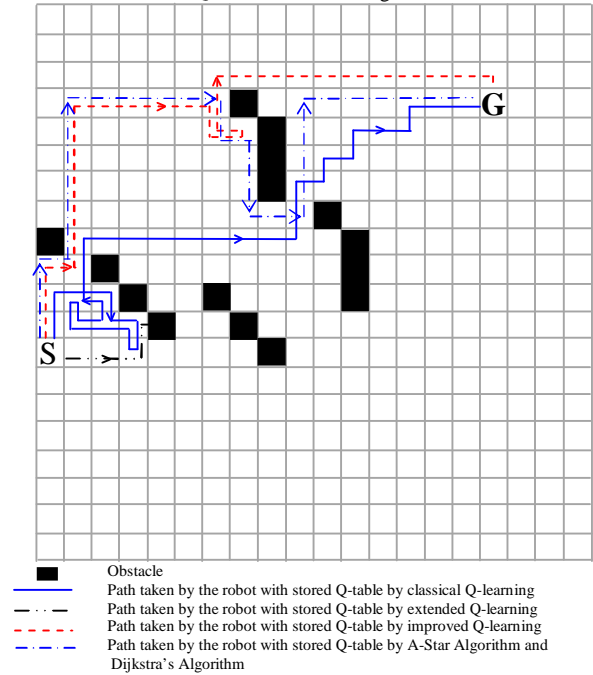


Fig 4: World map with obstacles. It shows the paths taken by the robot with the stored Q-table of different algorithms. The robot with the stored Q-table by the extended Q-learning fails to reach the goal.

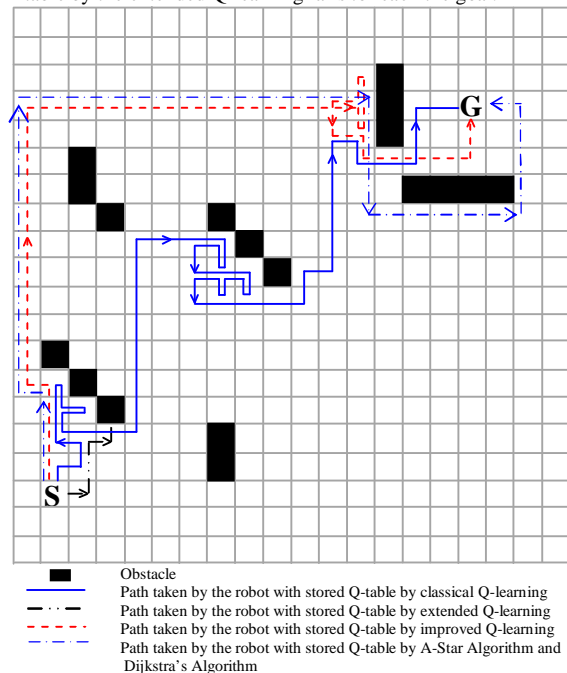


Fig 5: World map 5 with obstacles. It shows the path taken by the robot with the stored Q-table of different algorithms. Here the robot with the stored Q-table by the extended Q-learning fails to reach the goal.

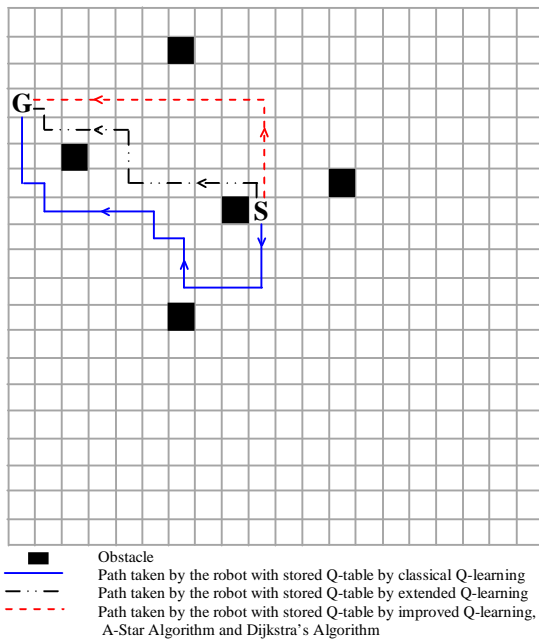


Fig 6: World map 6 with obstacles. It shows the paths taken by the robot with the stored Q-table of different algorithms.

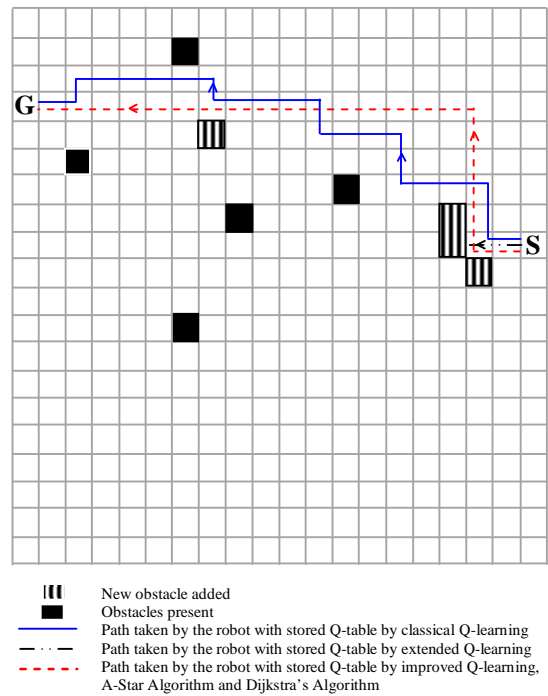


Fig 8: World map 8 with obstacles. Paths taken by the robot with stored Q-table by the different algorithms are shown in the fig. The robot with the stored Q-table by the extended Q-learning fails to reach the goal.

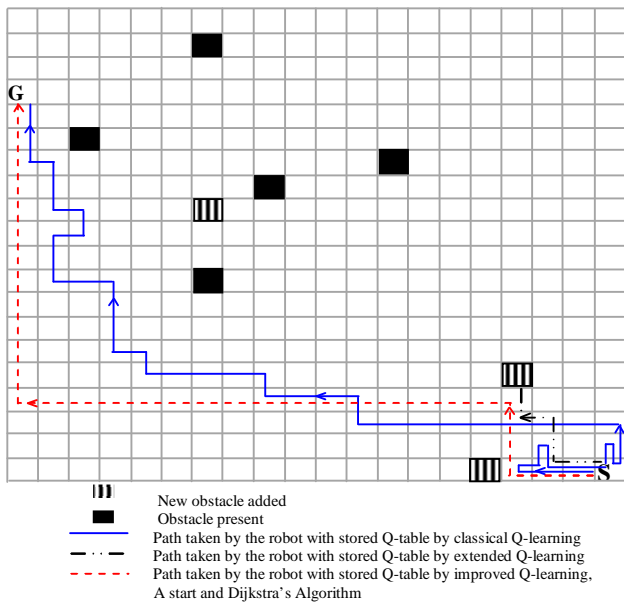


Fig 7: World map 7 with obstacles. Paths taken by the robot with stored Q-table by the different algorithms are shown in the fig. Here the robot with the stored Q-table by the extended Q-learning fails to reach the goal.

Table 1: Comparison of time taken by the robot and number of 90° turn required by the robot

World map	Planning time taken in seconds					No. of 90° turns				
	<i>IQL</i>	<i>EQL</i>	<i>CQL</i>	<i>A-star</i>	<i>Dijkstra</i>	<i>IQL</i>	<i>EQL</i>	<i>CQL</i>	<i>A-star</i>	<i>Dijkstra</i>
Fig. 1	22.25	24.82	26.47	40.25	88.25	2	15	16	2	2
Fig. 2	23.07	-	30.76	42.05	91.24	6	-	20	6	6
Fig. 3	13.78	14.17	16.70	27.65	52.60	4	4	10	4	4
Fig. 4	23.40	-	32.95	47.90	93.86	11	-	22	10	10
Fig. 5	31.48	-	47.08	64.01	121.35	15	-	38	8	8
Fig. 6	09.34	10.16	14.77	18.52	36.73	2	6	9	2	2
Fig. 7	24.55	-	37.51	49.47	122.31	5	-	29	5	5
Fig. 8	16.81	-	20.21	32.76	63.87	4	-	12	4	4

IQL Improved Q-learning
EQL Extended Q-learning
CQL Classical Q-learning
 - Goal cannot be reached

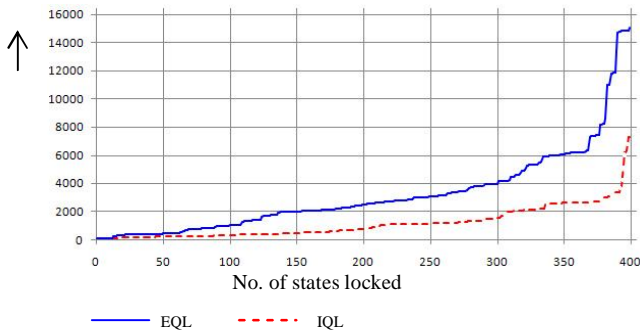


Fig. 9: Comparison between no. of iterations required by the improved Q-learning (IQL) algorithm and the extended Q-learning (EQL) algorithm to learn the world map (given in Fig. 2) without any obstacle.

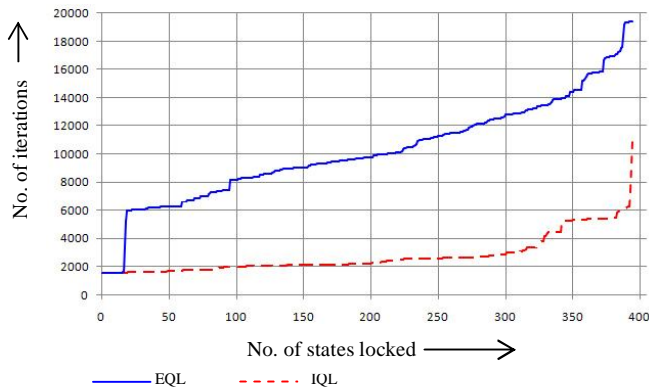


Fig. 10: Comparison between no of iteration required by the improved Q-learning (IQL) algorithm and the extended Q-learning (EQL) algorithm to learn the world map (given in Fig. 7) with 5 obstacles.

Fig. 9 and 10 provide information about the number of states locked for EQL and IQL. It is apparent from the figures that IQL takes relatively smaller number of iterations than EQL for having the same number of states locked. This justifies the significance of the two additional locking conditions in IQL when compared to EQL.

VII. EXPERIMENTS WITH KHEPERA II ROBOT

Khepera II (Fig. 11) is a miniature robot (diameter of 7 cm) equipped with 8 built-in infrared proximity sensors, and 2 relatively accurate encoders for the two motors. The range sensors are positioned at fixed angles and have limited range detection capabilities. The sensors are numbered between 0 and 7 with the leftmost sensor, designated by 0, and the rightmost by 7 (Fig. 12). The robot represents measured range data in the scale: $[0, 1023]$. When an obstacle is away from the sensor by more than 5cm, it is represented by zero. When an obstacle is approximately 2 cm away, it is represented by 1023. The onboard Microprocessor includes a flash memory of 512 KB, and a Motorola 68331, 25MHz processor. The Khepera model we used is a table-top robot, connected to a workstation through a wired serial link. This configuration allows an optional experimental configuration with everything at hand: the robot, the environment and the host computer.

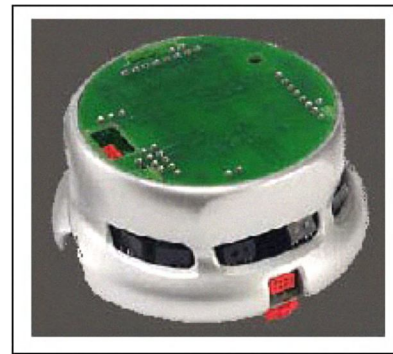


Figure 11: The Khepera II Robot

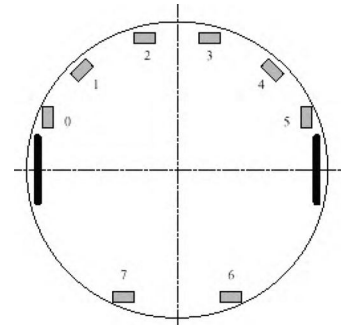


Fig. 12: Position of the sensors of Khepera II

Different experimental world maps have been developed to study the performance of the three different Q-learning and the corresponding path-planning algorithms. The starting and the goal states S and G are marked in all the experimental maps. The snapshots of each map after construction of the trajectory by the robot using distinctive colored lines for three algorithms CQL, EQL and IQL for all the experiments are given.

The first experiment is developed based on the world map shown in Fig. 13. The IQL and the EQL respectively takes 179 and 1710 iterations to learn the said environment. The CQL algorithm with random selection requires 3012 iterations for convergence. After the learning phase is over, the path planning algorithm is executed by keeping the Khepera at state no 48, facing left, in the experimental world map and is allowed to traverse to the goal using the stored Q-table by all the three algorithms. Path taken by the robot with the stored Q-table by the IQL, the CQL and the EQL are shown in Fig. 13 by red blue and green lines respectively.

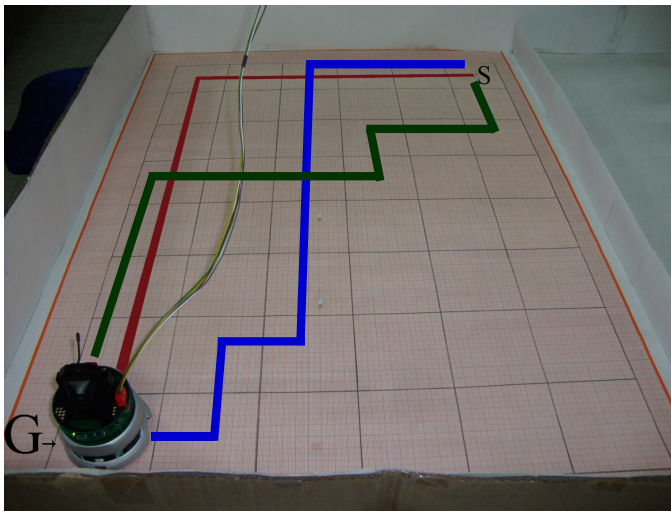


Fig.13: World map without any obstacle. Path planned by the robot:
 — Using the Q-table returned by the IQL.
 — Using the Q-table returned by the CQL.
 — Using the Q-table returned by the EQL.

The second experiment considers Q-learning in Fig. 14, and path-planning in the same world map after introducing six rectangular obstacles. The path taken by the robot with the stored Q-table by IQL, CQL and EQL are shown in Fig. 14 by distinctive colored lines. The robot with the stored Q-table by the EQL fails to reach the goal as indicated by black line in Fig 14. After reaching state no. 36, the best action stored in the Q-table is RIGHT action. In order to perform this action the robot turns right by 90° and checks for the presence of any obstacle in front of it with the help of sensor 2 and sensor 3. The robot finds an obstacle in state no. 35 and there is no alternative action stored in the Q-table. Therefore the robot with the stored Q-table obtained by the EQL fails to reach goal and stop at state no. 36.

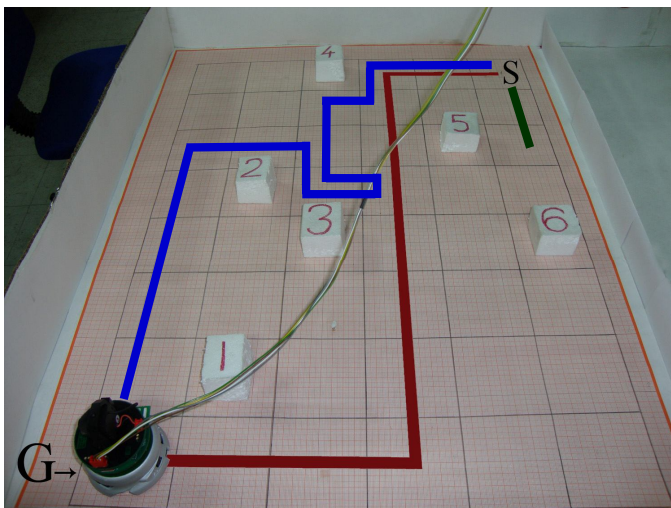


Fig. 14: World map with six obstacles added after the learning phase. Path planned by the robot:
 — Using the Q-table returned by the IQL.
 — Using the Q-table returned by the CQL.
 — Using the Q-table returned by the EQL (The robot fails to reach the goal. After reaching state no. 36 the next state is

state no. 35 but an obstacle is present in that state. The robot will stop at the state no 36).

Experiment 3 is concerned with learning and planning in the world map shown in Fig. 15 containing four obstacles, numbered 1 to 4. All the three algorithms are used to learn the movement steps from each grid in the map to its neighboring grid. The IQL takes 311 iterations to learn the said environment, while the EQL algorithm takes 1857 iterations for the same learning task. The CQL algorithm with random action selection takes 3060 iterations for convergence. After completion of the learning phase we kept the Khepera, facing left, in the same environment and the planning phase is executed with all three algorithms by using the stored Q-table obtained by the respective learning algorithm. The red, blue and green lines in Fig. 15 show the paths taken by the robot with the stored Q-table by IQL, CQL and EQL respectively.

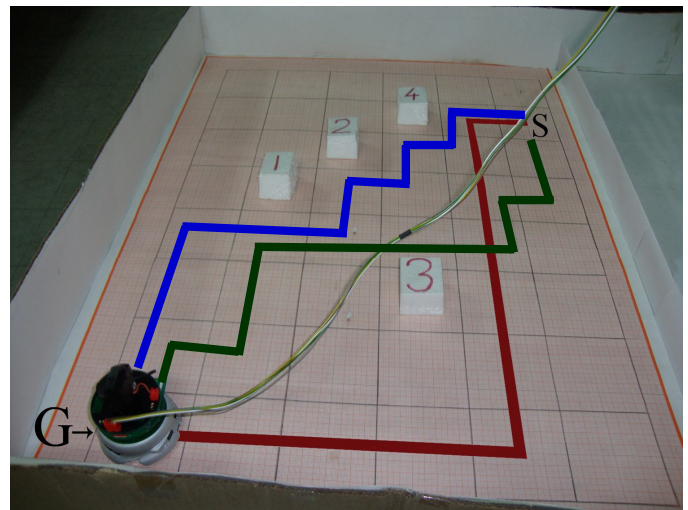


Fig. 15: World map with four obstacles. Path planned by the robot using:
 — Using the Q-table returned by the IQL.
 — Using the Q-table returned by the CQL.
 — Using the Q-table returned by the EQL

In experiment 4, we train the robot in the world map given in Fig. 16 with four obstacles numbered 1 to 4. After the training is over, we add two obstacles numbered 5 and 6 in the map, the planning cycle is executed in the modified world map given in Fig. 16. The path taken by the robot with the stored Q-table by IQL, CQL and EQL are given by red, blue and green lines respectively. The robot with the stored Q-table by the extended Q-learning fails to reach the goal as indicated by the incomplete green line segment in Fig. 16. The justification of the failure is given below. After reaching state no. 8, the robot determines the best action at this state, which is move left, but an obstacle is present at the next state. Therefore, the robot reports failure and stops in state no. 8.

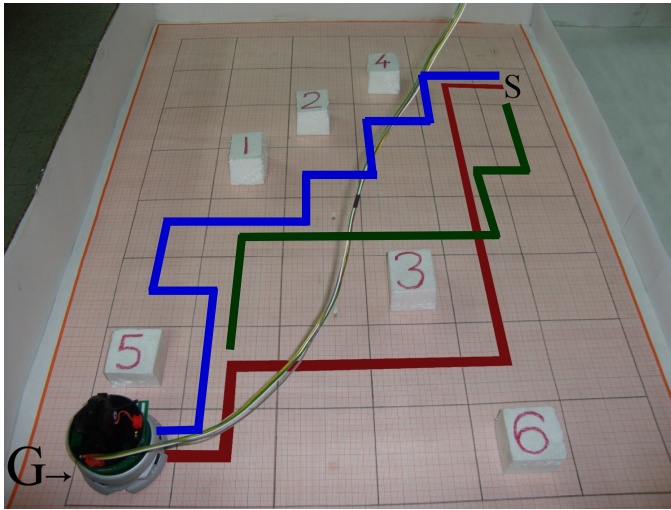


Fig. 16: World map with four obstacles. Path planned by the robot using:

- Using the Q-table returned by the IQL.
- Using the Q-table returned by the CQL.
- Using the Q-table returned by the EQL (When the robot reaches at state no. 8, the best action at state no. 8 is move left but an obstacle is present at state no. 7. So, the robot will stop at the state no.8).

Results of the experiments undertaken above are summarized in Table-2. The Table compares the relative performance of the IQL, EQL and CQL in path-planning. The metrics employed to compare the relative merits of the Q-learning algorithms in the planning phase are: 1) time taken to reach the goal, 2) the number of 90° turns involved in the path-planning, and 3) the number of states traversed during the planning phase.

It is apparent from Table-2 that for all the world maps shown in Fig. 13-16 the IQL outperforms the CQL and the EQL with respect to all the three metrics. The maximum noteworthy merit of the IQL over the others is the number of 90° turns, which is a bare minimum as evident from Table-2.

Table 2: Comparison of Time taken by the robot, no. of 90° turns required by the robot and no. of state traversed by the robot

Fig. No	Time Taken in sec			No. of 90° turn			No. of states traversed		
	IQL	EQL	CQL	IQL	EQL	CQL	IQL	EQL	CQL
13	40.73	48.44	47.40	1	5	4	12	12	12
14	43.74	-	71.17	2	-	9	12	-	16
15	39.62	48.40	48.40	2	7	7	11	11	11
16	43.72	-	59.52	4	-	10	11	-	13

IQL Improved Q-learning
EQL Extended Q-learning
CQL Classical Q-learning
- Goal cannot be reached

VIII. Conclusions

The paper proposed an alternative algorithm for deterministic Q-learning, presuming that the background knowledge about the distance from the current state to both the next state and the goal state are available. The proposed algorithm updates the entries of the Q-table only once unlike the classical Q-learning, where the entries in the Q-table were updated many

times until convergence was ensured. This results in a significant saving in time complexity of the order of mn in comparison to the classical Q-learning, where n and m are the number of states and number of actions at each state respectively.

Theorem 1 indicates the correctness in the steady-state values in the entries of the Q-table. Time-complexity analysis also reveals that the proposed algorithm saves a time-complexity of the order of mn , when compared to the classical Q-learning.

Experiments simulated on different experimental maze and on the Khepera platform confirms the better performance of the proposed algorithm in comparison to Classical and the extended Q-learning algorithms. The Q-table updated by the improved Q-learning, when used for path-planning application, outperforms both the classical and the extended Q-learning with respect to all the three metrics used in the experimental study. Most importantly, the 90° turnings required in the IQL is significantly reduced in comparison to the CQL and the EQL. Since the IQL outperforms CQL and EQL in all the three metrics as indicated in Table 1 and 2, it has a good potential in path-planning applications of mobile robots, particularly when obstacles are added in real time

REFERENCES

- [1] Dean, T., Basye, K. and Shewchuk, J. "Reinforcement learning for planning and Control". In: *Minton, S. (ed.) Machine Learning Methods for Planning and Scheduling*. Morgan Kaufmann 1993.
- [2] Bellman, R.E., *Dynamic programming*, Princeton, NJ: Princeton University Press, 1957.
- [3] Watkins, C. and Dayan, P., "Q-learning", *Machine Learning*, Vol. 8, pp. 279- 292, 1992
- [4] Konar, A., *Computational Intelligence: Principles, Techniques and Applications*, Springer-Verlag, 2005
- [5] Busoniu, L., Babushka, R., Schutter, B.De., Ernst, D., *Reinforcement Learning and Dynamic Programming Using Function Approximators*, CRC Press, Taylor & Francis group, Boca Raton, FL, 2010.
- [6] Chakraborty, J., Konar A., Jain, L.C., and Chakraborty, U., "Cooperative Multi-Robot Path Planning Using Differential Evolution" *Journal of Intelligent & Fuzzy Systems*, Vol. 20, Pp.13-27, 2009.
- [7] Gerke, M., and Hoyer, H., "Planning of Optimal paths for autonomous agents moving in inhomogeneous environments", in: *Proceedings of the 8th Int. Conf. on Advanced Robotics*, July 1997, pp.347-352.
- [8] Xiao, J., Michalewicz, Z., Zhang, L., and Trojanowski, K., "Adaptive Evolutionary Planner/ Navigator for Mobile robots", *IEEE Transactions on evolutionary Computation* 1 (1), April 1997.
- [9] Bien, Z., and Lee, J., "A Minimum-Time trajectory planning Method for Two Robots", *IEEE Trans on Robotics and Automation* 8(3), PP.443-450, JUNE 1992.
- [10] Moll, M., and Kavraki, L.E., "Path Planning for minimal Energy Curves of Constant Length", in: *Proceedings of the 2004 IEEE Int. Conf. on Robotics and Automation*, pp.2826-2831, April 2004.
- [11] Regele, R., and Levi, P., "Cooperative Multi-Robot Path Planning by Heuristic Priority Adjustment", in: *Proceedings of the IEEE/RJS Int Conf on Intelligent Robots and Systems*, 2006.
- [12] Yuan-Pao Hsu, Wei-Cheng Jiang, Hsin-Yi Lin, "A CMAC-Q-Learning Based Dyna Agent", in: *SICE Annual Conference*, 2008, pp. 2946 – 2950, The University Electro-Communications, Tokyo, Japan.
- [13] Yi Zhou and Meng Joo Er, "A Novel Q-Learning Approach with Continuous States and Actions", in: *16th IEEE International Conference on Control Applications Part of IEEE Multi-conference on Systems and Control*, Singapore, 1-3 October 2007

- [14] Kyungeun Cho, Yunsick Sung, Kyhyun Um, "A Production Technique for a Q-table with an Influence Map for Speeding up Q-learning", in : *International Conference on Intelligent Pervasive Computing*, 2007.
- [15] Deepshikha Pandey, Punit Pandey, "Approximate Q-Learning: An Introduction", in : *Second International Conference on Machine Learning and Computing*, 2010.
- [16] S. S. Masoumzadeh and G. Taghizadeh, K. Meshgi and S. Shiry, Deep Blue, "A Fuzzy Q-Learning Enhanced Active Queue Management Scheme", in : *International Conference on Adaptive and Intelligent Systems*, 2009.
- [17] Indrani Goswami (Chakraborty), Pradipta Kumar Das, Amit Konar, R. Janarthanan. "Extended Q-learning Algorithm for Path-Planning of a Mobile Robot", in: *Eighth International Conference on Simulated Evolution And Learning (SEAL-2010)*, Indian Institute of Technology Kanpur, India, December 2010.
- [18] L. A. Jeni, Z. Istenes P Korondi, H. Hashimoto Hierarchical "Reinforcement Learning for Robot Navigation using the Intelligent Space Concept", in *11th International Conference on Intelligent Engineering Systems - 29 June - 1 July 2007 - Budapest, Hungary*
- [19] Tom'as Mart'inez-Mar'in and Rafael Rodr'iguez, "Navigation of Autonomous Vehicles in Unknown Environments using Reinforcement Learning", in *2007 IEEE Intelligent Vehicles Symposium*, Istanbul, Turkey, June 13-15, 2007
- [20] Wooyoung Kwon, Il Hong Suh, Sanghoon Lee, Young-Jo Cho, "Fast Reinforcement Learning Using Stochastic Shortest Paths for a Mobile Robot" in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, San Diego, CA, USA, Oct 29 - Nov 2, 2007
- [21] L. A. Jeni , Z. Istenes , P'eter Szemes, H. Hashimoto, "Robot Navigation Framework Based on Reinforcement Learning for Intelligent Space", *2008 Conference on Human System Interaction*, Krakow, Poland, May 25-27, 2008
- [22] Roland Lang, Stefan Kohlhauser, Gerhard Zucker, and Tobias Deutsch, "Integrating Internal Performance Measures into the Decision Making Process of Autonomous Agents", in *2010 Conference on Human System Interaction*, Rzeszow, Poland, May 13-15, 2010.
- [23] G. Tesauro, "Extending Q-Learning to General Adaptive. Multi-Agent Systems", in *Advances in Neural Information Processing Systems*, volume 16, 2004.
- [24] Frazier, P. and W. B. Powell, "The Knowledge Gradient Policy for Offline Learning with Independent Normal Rewards", in *Proceedings of the 2007 IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning (ADPRL 2007)*.
- [25] Jung-Jun Park, Ji-Hun Kim, and Jae-Bok Song, "Path Planning for a Robot Manipulator based on Probabilistic Roadmap and Reinforcement Learning", in *International Journal of Control, Automation, and Systems*, vol. 5, no. 6, pp. 674-680, December 2007.
- [26] Shoufeng Lu, Ximin Liu, Shiqiang Dai, "Incremental Multistep Q-learning for Adaptive Traffic Signal Control Based on Delay Minimization Strategy", in *Proceedings of the 7th World Congress on Intelligent Control and Automation*, June 25 - 27, 2008, Chongqing, China.
- [27] Wei Chen, Jing Guo, Xiong Li, Jie Wang, "Hybrid Q-learning Algorithm About Cooperation in MAS", *2009 Chinese Control and Decision Conference (CCDC 2009)*.
- [28] E. Gomes and R. Kowalczyk, "Dynamic Analysis of Multiagent Q-learning with E-greedy Exploration," *Proceedings of the 26th International Conference on Machine Learning*, vol. 382, pp. 369-376, 2009.
- [29] Zhe Chen, and Robert C. Qiu, "Q-Learning Based Bidding Algorithm for Spectrum Auction in Cognitive Radio", in *Proceedings of IEEE SoutheastCon*, March, 2011.
- [30] A. Galindo-Serrano and L. Giupponi, "Distributed Q-learning for aggregated interference control in cognitive radio networks," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 4, pp. 1823 - 1834, 2010.
- [31] O. Alsaleh, B. Hamdaoui, and A. Fern, "Q-learning for opportunistic spectrum access," in *Proceedings of the 6th International Wireless Communications and Mobile Computing Conference, 2010*, pp. 220-224.
- [32] C. Wu, K. Chowdhury, and M. D. Felice, "Spectrum management of cognitive radio using multi-agent reinforcement learning," in *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, 2010.
- [33] R. C. Qiu, Z. Chen, N. Guo, Y. Song, P. Zhang, H. Li, and L. Lai, "Towards a real-time cognitive radio network testbed: architecture, hardware platform, and application to smart grid," in *Proceedings of the fifth IEEE Workshop on Networking Technologies for Software-Defined Radio and White Space*, June 2010.
- [34] Wang Yaping, Zhang Zheng, "A Method of Reinforcement Learning Based Automatic Traffic Signal Control", in *Third International Conference on Measuring Technology and Mechatronics Automation*, 2011
- [35] Hiroyuki Okamura and Tadashi Dohi, "Application of Reinforcement Learning to Software Rejuvenation", in *Tenth International Symposium on Autonomous Decentralized System 2011*
- [36] Sridhar Mahadevan, "Average Reward Reinforcement Learning: Foundations, Algorithms, and Empirical Results", *Machine Learning*, Special Issue on Reinforcement Learning (edited by Leslie Kaelbling), vol. 22, pp. 159-196, 1996.
- [37] L. Busoniu, R. Babuska, B.D. Schutter and D. Ernst, *Reinforcement Learning and Dynamic Programming using Function Approximators*, CRC Press, Florida, 2010.
- [38] Y. Sakaguchi and M. Takano, "Reliability of internal prediction/estimation and its application. I. Adaptive action selection reflecting reliability of value function, Neural Networks", vol. 17, pp. 935-952, 2004.
- [39] Sutton, R. S. and Barto, A. G., *Reinforcement Learning*, MIT Press, Boston, MA, 1998.
- [40] http://www.computationalintelligence.net/main/main_page.html



Amit Konar (SM'10) received the B.E. degree from Bengal Engineering and Science University (B.E. College), Howrah, India, in 1983 and the M.E. Tel E. M. Phil., and Ph.D. (Engineering) degrees from Jadavpur University, Calcutta-700032, India, in 1985, 1988, and 1994, respectively. In 2006, he was a Visiting Professor with the University of Missouri, St. Louis. He is currently a Professor with the Department of Electronics and Telecommunication Engineering (ETCE), Jadavpur University, where he is the Founding Coordinator of the M.Tech. program on intelligent automation and robotics. He has supervised fifteen Ph.D. theses. He has over 250 publications in international journal and conference proceedings.

He is the author of eight books, including two popular texts *Artificial Intelligence and Soft Computing* (CRC Press, 2000) and *Computational Intelligence: Principles, Techniques and Applications* (Springer, 2005). He serves as the Associate Editor of *IEEE Transactions on Systems, Man and Cybernetics, Part-A*, and *IEEE Transactions on Fuzzy Systems*. His research areas include the study of computational intelligence algorithms and their applications to the various domains of electrical engineering and computer science. Specifically, he worked on fuzzy sets and logic, neurocomputing, evolutionary algorithms, Dempster-Shafer theory, and Kalman filtering, and applied the principles of computational intelligence in image understanding, VLSI design, mobile robotics, pattern recognition, brain-computer interfacing and computational biology. He was the recipient of All India Council for

Technical Education (AICTE)-accredited 1997–2000 Career Award for Young Teachers for his significant contribution in teaching and research.



Indrani Goswami (Chakraborty) received her B.E. degree in Electrical engineering, M.E. degree (with specialization in Control Engineering) in Electronics and Communication Engineering and Ph.D. degree in Cognitive Robotics all from Jadavpur University in 1993, 1988 and 2012 respectively. She was a visiting Research scientist in Robotics Lab, MIE University, Nyoga, Japan. Indrani has been currently teaching in Calcutta Institute of Technology for the last 2 years.

Her current research interest includes Machine Intelligence and Robotics, Industrial Control, Cognitive science, Brain-Computer Interfacing and Evolutionary Algorithms. She has published a number of interesting papers in top international journals and conference proceedings.

Soft Computing (IJASIS) and serves on editorial boards for a number of prestigious journals as well as on International Programme Committee (IPC) for several international conferences. He received a prestigious Commonwealth Fellowship for pursuing his Doctorate in Applied Non-Linear Mathematics, which he earned from the University of York in 1996. He holds BSc (Hons.), MSc, and MPhil (with Distinction) from the MDS University of Ajmer, India. Prior to joining Liverpool Hope, Prof. Nagar was with the Department of Mathematical Sciences, and later at the Department of Systems Engineering, at Brunel University, London.



Sapam Jitu Singh received B.E. (Computer Technology) in 2000 from Nagpur University, Maharashtra, India and M.E. in Electronics & Tele-communication Engineering (Computer Engineering specialization) in 2011 from Jadavpur University, West Bengal, India. Since 2003, he has been working as Lecturer with Department of Computer Science and Engineering, Manipur Institute of Technology, Manipur, India. His research interests include evolutionary computing, digital image processing and robotics.



Lakhmi C. Jain is currently a Professor of knowledge-based engineering and the Director/Founder of the Knowledge-Based Intelligent Engineering Systems Centre, University of South Australia, Adelaide, Australia. His interests focus on artificial intelligence paradigms and their applications in complex systems, art-science fusion, e-education, e-healthcare, robotics, unmanned air vehicles, and intelligent agents. Prof. Jain is a Fellow of the Institution of Engineers Australia.



Atulya K. Nagar holds the Foundation Chair, as Professor of Computer and Mathematical Sciences, at Liverpool Hope University and is Head of the Department of Mathematics and Computer Science. A mathematician by training, Professor Nagar possesses multi-disciplinary expertise in Natural Computing, Bioinformatics, Operations Research and Systems Engineering. He has an extensive background and experience of working in Universities in the UK and India. He has been an expert reviewer for the Biotechnology and Biological Sciences Research Council (BBSRC) grants peer-review committee for Bioinformatics Panel and serves on Peer-Review College of the Arts and Humanities Research Council (AHRC) as a scientific expert member.

He has coedited volumes on Intelligent Systems, and Applied Mathematics; he is the Editor-in-Chief of the International Journal of Artificial Intelligence and

