

# Q Learning Algorithm

---

## AIM

---

To develop and evaluate the Q-learning algorithm's performance in navigating the environment and achieving the desired goal.

## PROBLEM STATEMENT

---

The goal of this project is to implement a Q-learning algorithm that enables an agent to learn optimal actions in a dynamic environment to maximize cumulative rewards.

## Q LEARNING ALGORITHM

---

### Step 1:

Initialize the Q-table with zeros for all state-action pairs based on the environment's observation and action space.

### Step 2:

Define the action selection method using an epsilon-greedy strategy to balance exploration and exploitation.

### Step 3:

Create decay schedules for the learning rate ( $\alpha$ ) and epsilon to progressively reduce their values over episodes.

### Step 4:

Loop through a specified number of episodes, resetting the environment at the start of each episode.

### Step 5:

Within each episode, continue until the episode is done, selecting actions based on the current state and the epsilon value.

### Step 6:

Execute the chosen action to obtain the next state and reward, and compute the temporal difference (TD) target.

### Step 7:

Update the Q-value for the current state-action pair using the TD error and the learning rate for that episode.

### Step 8:

Track the Q-values, value function, and policy after each episode for analysis and evaluation.

## Q LEARNING FUNCTION

---

**Name: Meetha Prabhu**

**Register Number: 212222240065**

**Importing necessary packages**

```
pip install git+https://github.com/mimoralea/gym-walk#egg=gym-walk
```

```
import warnings ; warnings.filterwarnings('ignore')
```

```
import itertools
import gym, gym_walk
import numpy as np
from tabulate import tabulate
from pprint import pprint
from tqdm import tqdm_notebook as tqdm
```

```
from itertools import cycle, count
```

```
import random
import matplotlib
import matplotlib.pyplot as plt
import matplotlib.pyplot as pylab
SEEDS = (12, 34, 56, 78, 90)
```

```
%matplotlib inline
```

**Function to plot the graph**

```
plt.style.use('fivethirtyeight')
params = {
    'figure.figsize': (5, 5),
    'font.size': 24,
    'legend.fontsize': 20,
    'axes.titlesize': 28,
    'axes.labelsize': 24,
    'xtick.labelsize': 20,
    'ytick.labelsize': 20
}
pylab.rcParams.update(params)
np.set_printoptions(suppress=True)
```

## Monte-Carlo Functions (Value\_iteration and Printing Functions)

```
def value_iteration(P, gamma=1.0, theta=1e-10):
    V = np.zeros(len(P), dtype=np.float64)
    while True:
        Q = np.zeros((len(P), len(P[0])), dtype=np.float64)
        for s in range(len(P)):
            for a in range(len(P[s])):
                for prob, next_state, reward, done in P[s][a]:
                    Q[s][a] += prob * (reward + gamma * V[next_state] * (not
done))
            if np.max(np.abs(V - np.max(Q, axis=1))) < theta:
                break
        V = np.max(Q, axis=1)
    pi = lambda s: {s:a for s, a in enumerate(np.argmax(Q, axis=1))}[s]
    return Q, V, pi

def print_policy(pi, P, action_symbols=('<', 'v', '>', '^'), n_cols=4,
title='Policy:'):
    print(title)
    arrs = {k:v for k,v in enumerate(action_symbols)}
    for s in range(len(P)):
        a = pi(s)
        print("| ", end="")
        if np.all([done for action in P[s].values() for _, _, _, done in
action]):
            print("".rjust(9), end=" ")
        else:
            print(str(s).zfill(2), arrs[a].rjust(6), end=" ")
        if (s + 1) % n_cols == 0: print("|")

def print_state_value_function(V, P, n_cols=4, prec=3, title='State-value
function:'):
    print(title)
    for s in range(len(P)):
        v = V[s]
        print("| ", end="")
```

```

        if np.all([done for action in P[s].values() for _, _, done in
action])):
            print("".rjust(9), end=" ")
        else:
            print(str(s).zfill(2), '{}'.format(np.round(v, prec)).rjust(6),
end=" ")
            if (s + 1) % n_cols == 0: print("|")

def print_action_value_function(Q,
                                optimal_Q=None,
                                action_symbols=('<', '>'),
                                prec=3,
                                title='Action-value function:'):
    vf_types=('',) if optimal_Q is None else ('', '*', 'err')
    headers = ['s',] + [' '.join(i) for i in
list(itertools.product(vf_types, action_symbols))]
    print(title)
    states = np.arange(len(Q))[..., np.newaxis]
    arr = np.hstack((states, np.round(Q, prec)))
    if not (optimal_Q is None):
        arr = np.hstack((arr, np.round(optimal_Q, prec), np.round(optimal_Q-
Q, prec)))
    print(tabulate(arr, headers, tablefmt="fancy_grid"))

def get_policy_metrics(env, gamma, pi, goal_state, optimal_Q,
                      n_episodes=100, max_steps=200):
    random.seed(123); np.random.seed(123) ; env.seed(123)
    reached_goal, episode_reward, episode_regret = [], [], []
    for _ in range(n_episodes):
        state, done, steps = env.reset(), False, 0
        episode_reward.append(0.0)
        episode_regret.append(0.0)
        while not done and steps < max_steps:
            action = pi(state)
            regret = np.max(optimal_Q[state]) - optimal_Q[state][action]
            episode_regret[-1] += regret

            state, reward, done, _ = env.step(action)
            episode_reward[-1] += (gamma**steps * reward)

            steps += 1

        reached_goal.append(state == goal_state)
    results = np.array((np.sum(reached_goal)/len(reached_goal)*100,
                        np.mean(episode_reward),
                        np.mean(episode_regret)))
    return results

```

## Metrics Functions

```
def get_metrics_from_tracks(env, gamma, goal_state, optimal_Q, pi, track
```

```

def get_metrics_from_tracks(env, gamma, goal_state, optimal_Q, pi_track,
coverage=0.1):
    total_samples = len(pi_track)
    n_samples = int(total_samples * coverage)
    samples_e = np.linspace(0, total_samples, n_samples, endpoint=True,
dtype=np.int)
    metrics = []
    for e, pi in enumerate(tqdm(pi_track)):
        if e in samples_e:
            metrics.append(get_policy_metrics(
                env,
                gamma=gamma,
                pi=lambda s: pi[s],
                goal_state=goal_state,
                optimal_Q=optimal_Q))
        else:
            metrics.append(metrics[-1])
    metrics = np.array(metrics)
    success_rate_ma, mean_return_ma, mean_regret_ma =
np.apply_along_axis(moving_average, axis=0, arr=metrics).T
    return success_rate_ma, mean_return_ma, mean_regret_ma

def rmse(x, y, dp=4):
    return np.round(np.sqrt(np.mean((x - y)**2)), dp)

def moving_average(a, n=100) :
    ret = np.cumsum(a, dtype=float)
    ret[n:] = ret[n:] - ret[:-n]
    return ret[n - 1:] / n

```

## Plotting Value Functions

```

def plot_value_function(title, V_track, V_true=None, log=False,
limit_value=0.05, limit_items=5):
    np.random.seed(123)
    per_col = 25
    linecycler = cycle(["-", "--", ":", "-."])
    legends = []

    valid_values = np.argwhere(V_track[-1] > limit_value).squeeze()
    items_idx = np.random.choice(valid_values,
                                min(len(valid_values), limit_items),
                                replace=False)

    # draw the true values first
    if V_true is not None:
        for i, state in enumerate(V_track.T):
            if i not in items_idx:
                continue
            if state[-1] < limit_value:
                continue

```

```

        label = 'v*({})'.format(i)
        plt.axhline(y=V_true[i], color='k', linestyle='--', linewidth=1)
        plt.text(int(len(V_track)*1.02), V_true[i]+.01, label)

# then the estimates
for i, state in enumerate(V_track.T):
    if i not in items_idx:
        continue
    if state[-1] < limit_value:
        continue
    line_type = next(linecycler)
    label = 'V({})'.format(i)
    p, = plt.plot(state, line_type, label=label, linewidth=3)
    legends.append(p)

legends.reverse()

ls = []
for loc, idx in enumerate(range(0, len(legends), per_col)):
    subset = legends[idx:idx+per_col]
    l = plt.legend(subset, [p.get_label() for p in subset],
                  loc='center right', bbox_to_anchor=(1.25, 0.5))
    ls.append(l)
[plt.gca().add_artist(l) for l in ls[:-1]]
if log: plt.xscale('log')
plt.title(title)
plt.ylabel('State-value function')
plt.xlabel('Episodes (log scale)' if log else 'Episodes')
plt.show()

```

## Decay\_schedule Function

```

def decay_schedule(init_value, min_value, decay_ratio, max_steps,
log_start=-2, log_base=10):
    decay_steps = int(max_steps * decay_ratio)
    rem_steps = max_steps - decay_steps
    values = np.logspace(log_start, 0, decay_steps, base=log_base,
endpoint=True)[::-1]
    values = (values - values.min()) / (values.max() - values.min())
    values = (init_value - min_value) * values + min_value
    values = np.pad(values, (0, rem_steps), 'edge')
    return values

```

## Slippery walk Seven

```

env = gym.make('SlipperyWalkSeven-v0')
init_state = env.reset()
goal state = 8

```

```

gamma = 0.99
n_episodes = 3000
P = env.env.P
n_cols, svf_prec, err_prec, avf_prec=9, 4, 2, 3
action_symbols=('<', '>')
limit_items, limit_value = 5, 0.0
cu_limit_items, cu_limit_value, cu_episodes = 10, 0.0, 100

```

## Alpha and Epsilon Functions

```

plt.plot(decay_schedule(0.5, 0.01, 0.5, n_episodes),
         '-', linewidth=2,
         label='Alpha schedule')
plt.plot(decay_schedule(1.0, 0.1, 0.9, n_episodes),
         ':', linewidth=2,
         label='Epsilon schedule')
plt.legend(loc=1, ncol=1)

plt.title('Alpha and epsilon schedules')
plt.xlabel('Episodes')
plt.ylabel('Hyperparameter values')
plt.xticks(rotation=45)

plt.show()

```

## Optimal value Function and policy

```

optimal_Q, optimal_V, optimal_pi = value_iteration(P, gamma=gamma)
print_state_value_function(optimal_V, P, n_cols=n_cols, prec=svf_prec,
                           title='Optimal state-value function:')
print()

print_action_value_function(optimal_Q,
                            None,
                            action_symbols=action_symbols,
                            prec=avf_prec,
                            title='Optimal action-value function:')

print()
print_policy(optimal_pi, P, action_symbols=action_symbols, n_cols=n_cols)
success_rate_op, mean_return_op, mean_regret_op = get_policy_metrics(
    env, gamma=gamma, pi=optimal_pi, goal_state=goal_state,
    optimal_Q=optimal_Q)
print('Reaches goal {:.2f}%. Obtains an average return of {:.4f}. Regret of
{:.4f}'.format(
    success_rate_op, mean_return_op, mean_regret_op))

```

## First visit Monte-Carlo Function

```
def generate_trajectory(select_action, Q, epsilon, env, max_steps=200):
    done, trajectory = False, []
    while not done:
        state = env.reset()
        for t in count():
            action = select_action(state, Q, epsilon)
            next_state, reward, done, _ = env.step(action)
            experience = (state, action, reward, next_state, done)
            trajectory.append(experience)
            if done:
                break
            if t >= max_steps - 1:
                trajectory = []
                break
            state = next_state
    return np.array(trajectory, object)

def mc_control(env,
               gamma=1.0,
               init_alpha=0.5,
               min_alpha=0.01,
               alpha_decay_ratio=0.5,
               init_epsilon=1.0,
               min_epsilon=0.1,
               epsilon_decay_ratio=0.9,
               n_episodes=3000,
               max_steps=200,
               first_visit=True):
    nS, nA = env.observation_space.n, env.action_space.n
    discounts = np.logspace(0,
                             max_steps,
                             num=max_steps,
                             base=gamma,
                             endpoint=False)

    alphas = decay_schedule(init_alpha,
                             min_alpha,
                             alpha_decay_ratio,
                             n_episodes)
    epsilons = decay_schedule(init_epsilon,
                              min_epsilon,
                              epsilon_decay_ratio,
                              n_episodes)

    pi_track = []
    Q = np.zeros((nS, nA), dtype=np.float64)
    Q_track = np.zeros((n_episodes, nS, nA), dtype=np.float64)
    select_action = lambda state, Q, epsilon: np.argmax(Q[state]) \
        if np.random.random() > epsilon \
        else np.random.randint(len(Q[state]))
```



```

for e in tqdm(range(n_episodes), leave=False):

    trajectory = generate_trajectory(select_action,
                                    Q,
                                    epsilons[e],
                                    env,
                                    max_steps)
    visited = np.zeros((nS, nA), dtype=bool)
    for t, (state, action, reward, _, _) in enumerate(trajectory):
        if visited[state][action] and first_visit:
            continue
        visited[state][action] = True

        n_steps = len(trajectory[t:])
        G = np.sum(discounts[:n_steps] * trajectory[t:, 2])
        Q[state][action] = Q[state][action] + alphas[e] * (G - Q[state]
[action])

    Q_track[e] = Q
    pi_track.append(np.argmax(Q, axis=1))

V = np.max(Q, axis=1)
pi = lambda s: {s:a for s, a in enumerate(np.argmax(Q, axis=1))}[s]
return Q, V, pi, Q_track, pi_track

Q_mcs, V_mcs, Q_track_mcs = [], [], []
for seed in tqdm(SEEDS, desc='All seeds', leave=True):
    random.seed(seed); np.random.seed(seed) ; env.seed(seed)
    Q_mc, V_mc, pi_mc, Q_track_mc, pi_track_mc = mc_control(env,
gamma=gamma, n_episodes=n_episodes)
    Q_mcs.append(Q_mc) ; V_mcs.append(V_mc) ; Q_track_mcs.append(Q_track_mc)
Q_mc, V_mc, Q_track_mc = np.mean(Q_mcs, axis=0), np.mean(V_mcs, axis=0),
np.mean(Q_track_mcs, axis=0)
del Q_mcs ; del V_mcs ; del Q_track_mcs

```

## Printing the value Functions

```

Q_mcs, V_mcs, Q_track_mcs = [], [], []
for seed in tqdm(SEEDS, desc='All seeds', leave=True):
    random.seed(seed); np.random.seed(seed) ; env.seed(seed)
    Q_mc, V_mc, pi_mc, Q_track_mc, pi_track_mc = mc_control(env,
gamma=gamma, n_episodes=n_episodes)
    Q_mcs.append(Q_mc) ; V_mcs.append(V_mc) ; Q_track_mcs.append(Q_track_mc)
Q_mc, V_mc, Q_track_mc = np.mean(Q_mcs, axis=0), np.mean(V_mcs, axis=0),
np.mean(Q_track_mcs, axis=0)
del Q_mcs ; del V_mcs ; del Q_track_mcs

```

## Q-learning

```

from tqdm import tqdm_notebook as tqdm
def q_learning(env,
               gamma=1.0,
               init_alpha=0.5,
               min_alpha=0.01,
               alpha_decay_ratio=0.5,
               init_epsilon=1.0,
               min_epsilon=0.1,
               epsilon_decay_ratio=0.9,
               n_episodes=3000):
    nS, nA = env.observation_space.n, env.action_space.n
    pi_track = []
    Q = np.zeros((nS, nA), dtype=np.float64)
    Q_track = np.zeros((n_episodes, nS, nA), dtype=np.float64)
    select_action = lambda state, Q, epsilon: np.argmax(Q[state]) \
        if np.random.random() > epsilon \
        else np.random.randint(len(Q[state]))
    alphas = decay_schedule ( init_alpha, min_alpha, alpha_decay_ratio,
                             n_episodes)

    epsilons = decay_schedule(init_epsilon,
                              min_epsilon,
                              epsilon_decay_ratio,
                              n_episodes)
    for e in tqdm(range(n_episodes), leave=False): # using tqdm
        state, done = env.reset(), False
        while not done:
            action = select_action(state, Q, epsilons[e])
            next_state, reward, done, _ = env.step(action)
            td_target = reward + gamma * Q[next_state].max() * (not done)
            td_error = td_target - Q[state][action]
            Q[state][action] = Q[state][action] + alphas[e] * td_error
            state = next_state
        Q_track[e] = Q
        pi_track.append(np.argmax(Q, axis=1))
    V=np.max(Q, axis=1)
    pi=lambda s: {s:a for s, a in enumerate(np.argmax(Q, axis=1))}[s]

    # Write your code here

    return Q, V, pi, Q_track, pi_track

from tqdm import tqdm_notebook as tqdm
Q_qls, V_qls, Q_track_qls = [], [], []
for seed in tqdm(SEEDS, desc='All seeds', leave=True):
    random.seed(seed); np.random.seed(seed) ; env.seed(seed)
    Q_ql, V_ql, pi_ql, Q_track_ql, pi_track_ql = q_learning(env,
gamma=gamma, n_episodes=n_episodes)
    Q_qls.append(Q_ql) ; V_qls.append(V_ql) ; Q_track_qls.append(Q_track_ql)
Q_ql = np.mean(Q_qls, axis=0)
V_ql = np.mean(V_qls, axis=0)
Q_track_ql = np.mean(Q_track_qls, axis=0)

```

```
Q_track_ql = np.mean(Q_track_qls, axis=0),
del Q_qls ; del V_qls ; del Q_track_qls
```

## Print the policy

```
print('Name:Meetha Prabhu      Register Number: 21222240065      ')
print_state_value_function(V_ql, P, n_cols=n_cols,
                           prec=svf_prec, title='State-value function found
by Q-learning:')
print_state_value_function(optimal_V, P, n_cols=n_cols,
                           prec=svf_prec, title='Optimal state-value
function:')
print_state_value_function(V_ql - optimal_V, P, n_cols=n_cols,
                           prec=err_prec, title='State-value function
errors:')
print('State-value function RMSE: {}'.format(rmse(V_ql, optimal_V)))
print()
print_action_value_function(Q_ql,
                            optimal_Q,
                            action_symbols=action_symbols,
                            prec=avf_prec,
                            title='Q-learning action-value function:')
print('Action-value function RMSE: {}'.format(rmse(Q_ql, optimal_Q)))
print()
print_policy(pi_ql, P, action_symbols=action_symbols, n_cols=n_cols)
success_rate_ql, mean_return_ql, mean_regret_ql = get_policy_metrics(
    env, gamma=gamma, pi=pi_ql, goal_state=goal_state, optimal_Q=optimal_Q)
print('Reaches goal {:.2f}%. Obtains an average return of {:.4f}. Regret of
{:.4f}'.format(
    success_rate_ql, mean_return_ql, mean_regret_ql))
```

## Plot for First-visit Monte-carlo

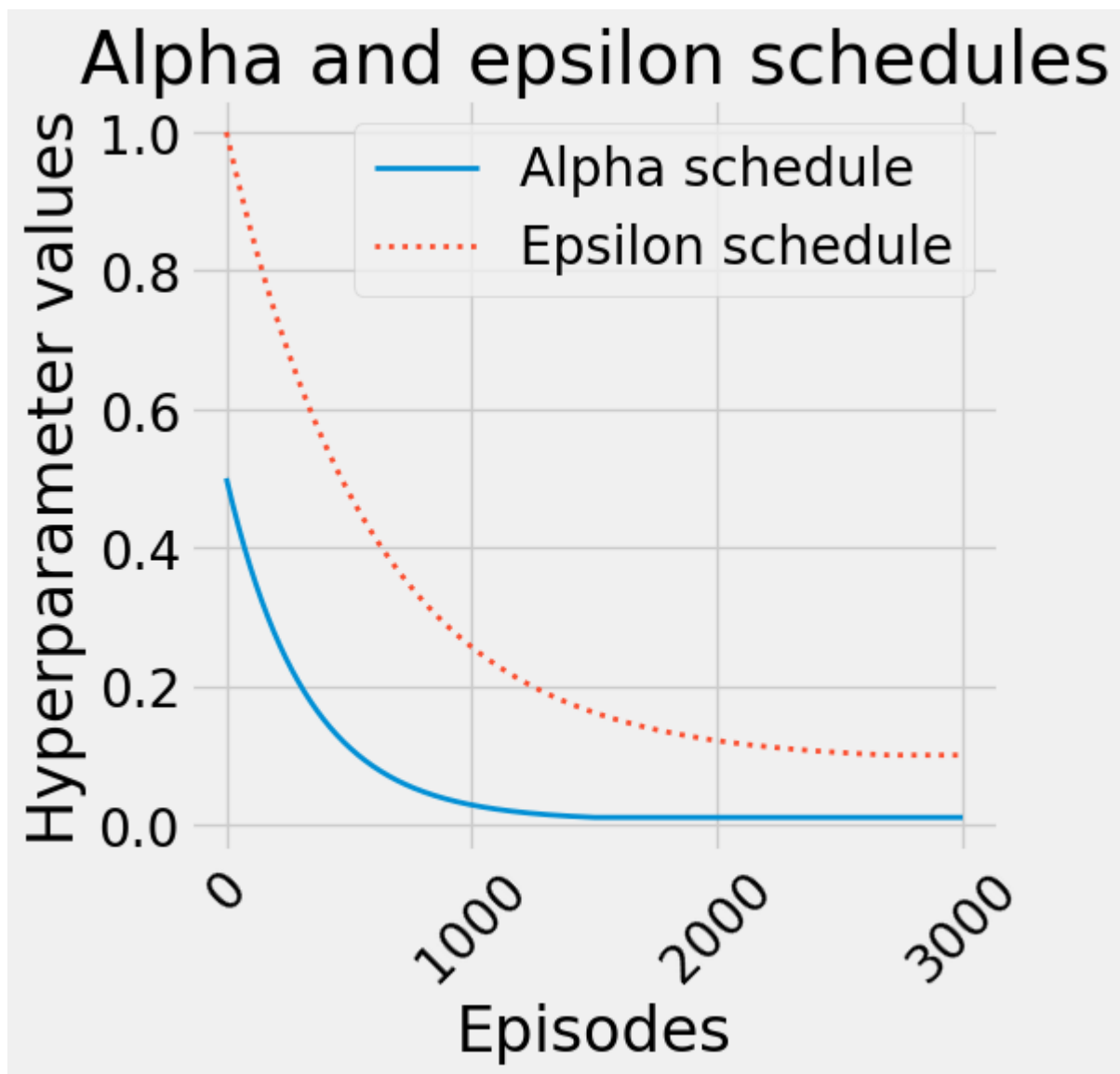
```
# First-visit Monte-carlo
plot_value_function(
    'FVMC estimates through time vs. true values      Name:Meetha Prabhu
Register Number: 21222240065      ',
    np.max(Q_track_mc, axis=2),
    optimal_V,
    limit_items=limit_items,
    limit_value=limit_value,
    log=False)
```

## Plot for First-visit Q-learning

# Q-learning

```
plot_value_function(  
    'Q-Learning estimates through time vs. true values', Name:Meetha  
    Prabhu, Register Number: 212222240065, '  
    np.max(Q_track_ql, axis=2),  
    optimal_V,  
    limit_items=limit_items,  
    limit_value=limit_value,  
    log=False)
```

## OUTPUT:



Optimal state-value function:

| 01 0.5637 | 02 0.763 | 03 0.8449 | 04 0.8892 | 05 0.922 | 06 0.9515 | 07 0.9806 |

Optimal action-value function:

s	<	>
0	0	0
1	0.312	0.564
2	0.67	0.763

3	0.803	0.845
4	0.864	0.889
5	0.901	0.922
6	0.932	0.952
7	0.961	0.981
8	0	0

Policy:

| 01 > | 02 > | 03 > | 04 > | 05 > | 06 > | 07 > |  
Reaches goal 96.00%. Obtains an average return of 0.8672. Regret of 0.0000

Name: Meetha Prabhu Register Number: 212222240065

State-value function found by FVMC:

| 01 0.502 | 02 0.7282 | 03 0.8219 | 04 0.8735 | 05 0.9146 | 06 0.9457 | 07 0.9797 |

Optimal state-value function:

| 01 0.5637 | 02 0.763 | 03 0.8449 | 04 0.8892 | 05 0.922 | 06 0.9515 | 07 0.9806 |

State-value function errors:

| 01 -0.06 | 02 -0.03 | 03 -0.02 | 04 -0.02 | 05 -0.01 | 06 -0.01 | 07 -0.0 |

State-value function RMSE: 0.0256

FVMC action-value function:

s	<	>	* <	* >	err <	err >
0	0	0	0	0	0	0
1	0.175	0.502	0.312	0.564	0.137	0.062
2	0.557	0.728	0.67	0.763	0.114	0.035
3	0.735	0.822	0.803	0.845	0.068	0.023
4	0.84	0.874	0.864	0.889	0.024	0.016
5	0.889	0.915	0.901	0.922	0.013	0.007
6	0.918	0.946	0.932	0.952	0.014	0.006
7	0.955	0.98	0.961	0.981	0.006	0.001
8	0	0	0	0	0	0

Action-value function RMSE: 0.049

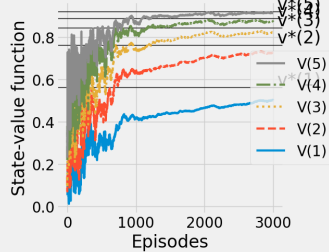
Policy:

| 01 > | 02 > | 03 > | 04 > | 05 > | 06 > | 07 > |  
Reaches goal 96.00%. Obtains an average return of 0.8672. Regret of 0.0000

FVMC estimates through time vs. true values

Name: Meetha Prabhu

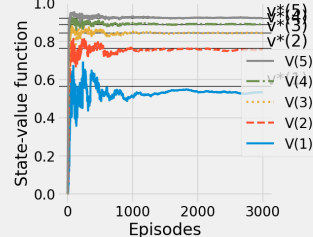
Register Number: 212222240065



Q-Learning estimates through time vs. true values

Name: Meetha Prabhu

Register Number: 212222240065



RESULT.

## RESULT:

Thus the python program to implement Q-learning is implemented successfully