

```

function Q=MonteCarloPolicyIteration(L, M, options)
    nstates = 3^9;    % 状態数
    nactions = 9;     % 行動数
    T = 5;           % 最大ステップ数

    % Q関数の初期化
    Q=zeros(nstates, nactions);
    Q=sparse(Q);

    % 政策反復
    for l=1:L
        visits = ones(nstates, nactions);    % (s, a) の出現回数
        results = zeros(M, 1);               % ゲームの結果
        rand('state', 1);                     % seedの初期化

        % エピソード
        for m=1:M
            state3 = zeros(1, 9);

            % ステップ
            for t=1:T

                % 状態のエンコード
                state = encode(state3);

                % 政策の生成
                policy = zeros(1, nactions);

                switch(options.pmode)
                    case 1 % greedy
                        [v, a] = max(Q(state, :));
                        policy(a) = 1;

                    case 2 % e-greedy
                        [v, a] = max(Q(state, :));
                        policy = ones(1, nactions)*options.epsilon/nactions;
                        policy(a) = 1-options.epsilon+options.epsilon/nactions;

                    case 3 % softmax
                        policy=exp(Q(state, :)/options.tau)/sum(exp(Q(state, :). /options.tau));
                end

                % 行動の選択および実行
                [action, reward, state3, fin] = action_train(policy, t, state3);

                % 状態, 行動, 報酬, 出現回数の更新
                states(m, t) = state;
                actions(m, t) = action;
                rewards(m, t) = reward;
                visits(state, action) = visits(state, action) + 1;

                % ゲーム終了
                if(fin>0)
                    results(m) = fin;
                end
            end
        end
    end
end

```

```

% 割引き報酬和の計算
drewards(m, t) = rewards(m, t);
for pstep=t-1:-1:1
    drewards(m, pstep) = options.gamma * drewards(m, pstep+1);
end
break;
end
end
end

% 状態行動価値関数の計算
Q=zeros(nstates, nactions);
Q=sparse(Q);
for m=1:M
    for t=1:size(states, 2)
        s = states(m, t);
        a = actions(m, t);
        if(s==0)
            break;
        end
        Q(s, a) = Q(s, a) + drewards(m, t);
    end
end

Q = Q ./ visits;

% 勝率の計算
rate(l) = size(find(results==2), 1) ./ M;

% 標準出力
fprintf(1, '%d Win=%d/%d, Draw=%d/%d, Lose=%d/%d\n', l, size(find(results==2), 1), M, size(find(
(results==3), 1), M, size(find(results==1), 1), M);
% fflush(stdout);
end

% グラフの出力
figure(1)
clf
% axes('FontSize', 15, 'LineWidth', 2.0);
games = M:M:M*L;
g=plot(games, rate);
set(g, 'LineWidth', 2);
g=xlabel('ゲーム数');
set(g, 'FontSize', 14);
g=ylabel('勝率');
set(g, 'FontSize', 14);
axis([M, M*L, 0.4, 1])

```