

# MATH4406 – Assignment 4

Patrick Laub (ID: 42051392)

September 26, 2014

## 1 Real-world motivation

A simple example of this type of control system is the thermostat for an air-conditioning unit. Once set to a desired temperature, the thermostat must enable and disable the air-conditioning over time to try to meet the set-point.

Consider the case where the thermostat is keeping an expensive server room at the correct temperature; if the temperature becomes too hot or too cold then the equipment could suffer costly damages, or operate inefficiently (hence the large negative rewards in state  $\pm 5$ ). Perhaps the other  $r(\cdot, \cdot)$  effects can be explained by additional costs accrued when the unit is swapped between heating and cooling mode.

When the air-conditioning is turned on ( $a = 1$ ) then the system is likely to cool down. Denote the random states  $X_i \in \{-5, \dots, 5\}$ , corresponding to system temperatures, then the previous statement means

$$\mathbb{P}(X_{n+1} > X_n | a = 1) \geq \mathbb{P}(X_{n+1} < X_n | a = 1)$$

and conversely the temperature is likely to drop when the air-condition is turned off ( $a = -1$ ):

$$\mathbb{P}(X_{n+1} < X_n | a = -1) \geq \mathbb{P}(X_{n+1} > X_n | a = -1).$$

## 2 Expected behaviour of optimal policies

Optimal policies will likely be symmetric about the state 0. For symmetric policies then the selected action in 0 will have no effect on overall system performance (as  $r(0, \cdot) = 0$ ).

When  $\lambda \downarrow 0$  then this means that future states become increasingly unimportant, and so the optimal policy should try to maximised the one-step expected payoff. If only looking at one-step ahead then a maximising policy should be (ignoring states  $\{-5, 0, 5\}$ )

$$d(s) = \begin{cases} -1, & s < 0 \\ 1, & s > 0 \end{cases}$$

However when  $\lambda \uparrow 1$  then the policy should be very conservative; the effects of the distant future are extremely important so the system should do everything it can to stick around  $s = 0$  (every  $X_t \in \{-5, 5\}$  will have a huge negative cost which should be avoided). A hypothesed optimal policy here is (ignoring states  $\{-5, 0, 5\}$ )

$$d(s) = \begin{cases} 1, & s < 0 \\ -1, & s > 0 \end{cases}$$

## 3 Optimal Policies: Discussion and Results

Brute-force enumeration, value iteration and policy iteration were used to solve the MDPs with  $\lambda \in \{0.01, 0.02, \dots, 0.99\}$ . Value iteration used a tolerance value of  $\epsilon = 10^{-9}$ . Each of the algorithms gave the same optimal policies (as expected). Fig. 1 visually<sup>1</sup> shows the optimal decision to take at every state  $s$  for every discounting factor  $\lambda$ . The policies with  $\lambda \approx 0$  and  $\lambda \approx 1$  match the expected/intuitive optimal policies outlined in Section 2.

---

<sup>1</sup>If a table is required, then the code in Section 4.4 constructs a table of the results in the matrix  $dRules$  which can easily be printed.

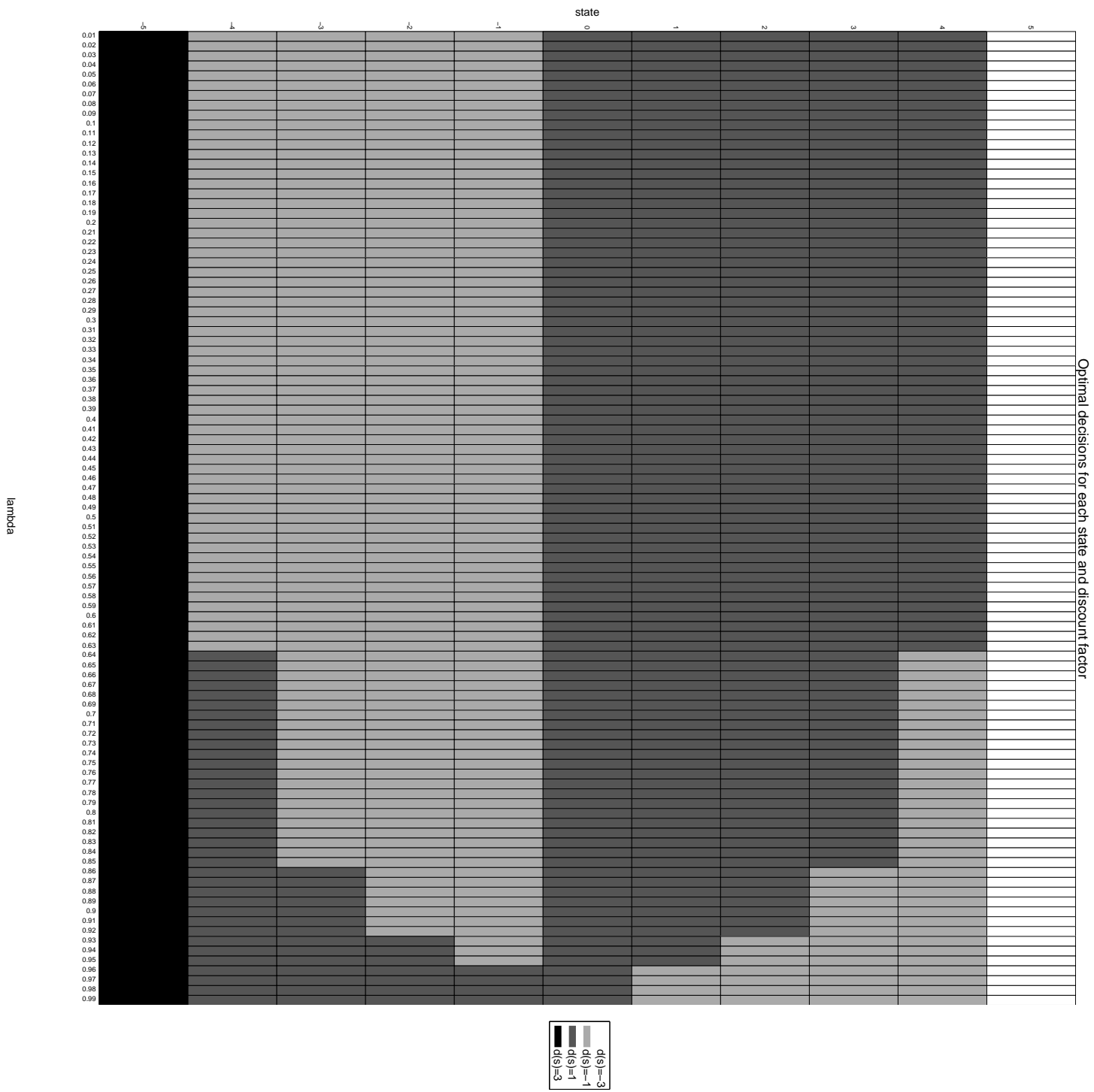


Figure 1: Optimal policies

## 4 Appendix: MATLAB Implementation

### 4.1 Brute force enumeration

```
1 function [v, d] = brute_force(lambda)
2     % Problem data.
3     S = (-5:5)'; numS = numel(S);
4
5     % Generate all decision rules.
6     allDecRules = allcomb(3, [-1,1], [-1,1], [-1,1], [-1,1], 1, ...
7         [-1,1], [-1,1], [-1,1], [-1,1], -3);
8
9     % Store the maximising value and decision rule so far.
10    bestV = -Inf .* ones(numS, 1);
11    bestD = NaN;
12
13    % Perform policy evaluation for every decision rule.
14    for ruleNum = 1:size(allDecRules, 1)
15        d = allDecRules(ruleNum, :)';
16
17        % Construct reward and transition vectors.
18        rd = S .* d;
19        Pd = zeros(numS);
20        Pd(1,1) = 0.5; Pd(1,2) = 0.5;
21        Pd(end,end) = 0.5; Pd(end,end-1) = 0.5;
22        for i = 2:(numS-1)
23            Pd(i,i-1) = 0.75 * (d(i) == -1) + 0.25 * (d(i) == 1);
24            Pd(i,i+1) = 0.25 * (d(i) == -1) + 0.75 * (d(i) == 1);
25        end
26
27        % Compute the value of the MDP.
28        v = (eye(numS) - lambda .* Pd) \ rd;
29
30        % Store if better (or as good as) the previous best.
31        if all((v - bestV) >= 0)
32            bestV = v;
33            bestD = d;
34        end
35    end
36
37    v = bestV; d = bestD;
38 end
```

## 4.2 Value iteration

```
1 function [v, d] = value_iteration(lambda)
2     % Problem data.
3     S = (-5:5)'; numS = numel(S);
4     eps = 1e-9;
5
6     % Start with any guess at the value.
7     v = zeros(numS, 1);
8     nextV = zeros(numS, 1);
9
10    % Construct decision rule at each step.
11    d = zeros(numS, 1); d(1) = 3; d(end) = -3; d(S == 0) = 1;
12
13    % Iterate over the value space until convergence (or give up).
14    for iters = 1:1e5
15        % For each state compute V^(n+1)(s) from V^n.
16        for i = 1:numS
17            % Handle edge cases separately to the middle states.
18            if i == 1
19                nextV(1) = S(1)*3 + lambda*(0.5*v(1) + 0.5*v(2));
20            elseif i == numS
21                nextV(end) = S(end)*-3 + lambda*(0.5*v(end) + 0.5*v(end-1));
22            else
23                up = S(i)*1 + lambda*(0.25*v(i-1) + 0.75*v(i+1));
24                down = S(i)*-1 + lambda*(0.75*v(i-1) + 0.25*v(i+1));
25
26                % Take the action which maximises this quantity. Note that
27                % when many maximisers exist, then this will select the
28                % larger one every time.
29                if up >= down
30                    nextV(i) = up; d(i) = 1;
31                else
32                    nextV(i) = down; d(i) = -1;
33                end
34            end
35        end
36
37        % Check for convergence.
38        if all(abs(nextV - v) < eps*(1-lambda)/(2*lambda))
39            v = nextV;
40            break;
41        end
42
43        v = nextV;
44    end
45
46    if iters == 1e5
47        warning('Did not converge!');
48    end
49 end
```

### 4.3 Policy iteration

```
1 function [v, d] = policy.iteration(lambda)
2   % Problem data.
3   S = (-5:5)'; numS = numel(S);
4
5   % Start with any guess decision rule.
6   d = ones(numS, 1); d(1) = 3; d(end) = -3;
7   nextD = d;
8
9   % Construct decision rule at each step.
10  for iters = 1:1e5
11    % Construct reward and transition vectors.
12    rd = S .* d;
13    Pd = zeros(numS);
14    Pd(1,1) = 0.5; Pd(1,2) = 0.5;
15    Pd(end,end) = 0.5; Pd(end,end-1) = 0.5;
16    for i = 2:(numS-1)
17      Pd(i,i-1) = 0.75 * (d(i) == -1) + 0.25 * (d(i) == 1);
18      Pd(i,i+1) = 0.25 * (d(i) == -1) + 0.75 * (d(i) == 1);
19    end
20
21    % Do policy evaluation.
22    v = (eye(numS) - lambda .* Pd) \ rd;
23
24    % Do policy improvement.
25    for i = 1:numS
26      if any(S(i) == [-5, 0, 5]), continue; end;
27
28      up = S(i)*1 + lambda*(0.25*v(i-1) + 0.75*v(i+1));
29      down = S(i)*-1 + lambda*(0.75*v(i-1) + 0.25*v(i+1));
30      if up >= down
31        nextD(i) = 1;
32      else
33        nextD(i) = -1;
34      end
35    end
36
37    % If the decision rule has converged then stop iterating.
38    if all(nextD == d)
39      break;
40    end
41
42    d = nextD;
43  end
44
45  if iters == 1e5
46    warning('Did not converge!');
47  end
48 end
```

## 4.4 Tests and plotting

```
1
2 dRules = [];
3
4 for lambda=0.01:0.01:0.99
5
6     [vb, db] = brute_force(lambda);
7     [vv, dv] = value_iteration(lambda);
8     [vp, dp] = policy_iteration(lambda);
9
10    if any(db ~= dv) || any(dv ~= dp)
11        error('Difference!!');
12    end
13
14    dRules = [dRules, db];
15 end
16
17 % Some grid code from: http://stackoverflow.com/questions/8711971
18 figure(1); clf;
19
20 % plot dummy objects
21 hold on;
22 for a=[-3, -1, 1, 3]
23     plot(1, 1, 'LineWidth', 10, 'Color', [(3-a)/6, (3-a)/6, (3-a)/6]);
24 end
25
26 % Draw large bounding box:
27 xstart = 0.01-0.01/2; ystart = -5-0.5;
28 xlen = 0.99; ylen = 11;
29
30 rectangle('position', [xstart, ystart, xlen, ylen])
31 dx = 0.01; dy = 1;
32 nx = floor(xlen/dx); ny = floor(ylen/dy);
33
34 for i = 1:nx
35     x = xstart + (i-1)*dx;
36     for j = 1:ny
37         y = ystart + (j-1)*dy;
38
39         a = dRules(j,i);
40         rectangle('position', [x, y, dx, dy], 'FaceColor', [(3-a)/6, (3-a)/6, (3-a)/6]);
41     end
42 end
43
44 axis([0.005, 0.994, -5.5, 5.499]);
45 set(gca, 'XTick', 0.01:0.01:0.99)
46 set(gca, 'YTick', -5:5)
47 xlabel('lambda'); ylabel('state');
48 title('Optimal decisions for each state and discount factor');
49 rotateXLabels( gca(), 90 )
50 legend('d(s)=-3', 'd(s)=-1', 'd(s)=1', 'd(s)=3')
```