Chap. 16

decision-theoretic agent = utility theory + prob. theory
    continuous measure of outcome quality

\* maximize expected utility

agent may not know current state

RESULT $(a)$ ≡ random variable
    values : possible outcome states

$$P(\text{RESULT}(a) = s' \mid a, \bar{e}) = \text{prob. of outcome } s' \text{ given } \bar{e}$$

$$= \sum_s P(\text{RESULT}(s,a) = s' \mid a) P(s_0 = s \mid \bar{e})$$

utility function    $U(s)$       desirability of state

expected utility    $EU(a \mid \bar{e}) = \sum_{s'} P(\text{RESULT}(a) = s' \mid a, e) U(s')$
         eqn (16.1)

principle of maximum expected utility (MEU)

$$\text{action} = \underset{a}{\text{argmax}} \; EU(a \mid \bar{e})$$

if utility captures performance measure

      then agent performs well

constraints on rational preferences

$A \succ B$      prefers A over B
$A \sim B$      indifferent between A + B
$A \succsim B$      indifferent or prefers A over B

consider set of outcomes for each action as a <u>lottery</u>
    L: outcomes $S_1, \ldots S_n$ with probs $P_1, \ldots P_n$

$$L = [P_1, S_1 \; ; \; P_2, S_2 \; ; \; \ldots \; P_n, S_n]$$

goal: understand how preferences between lotteries
     are related to    "      " underlying states

constraints of preference relation   (axioms of utility theory)

- Orderability    $(A \succ B) \lor (B \succ A) \lor (A \sim B)$
- Transitivity    $(A \succ B) \land (B \succ C) \Rightarrow (A \succ C)$
- Continuity    $A \succ B \succ C \Rightarrow \exists p \; [p, A \; ; \; 1-p, C] \sim B$
- Substitutability   $A \sim B \Rightarrow [p, A \; ; \; 1-p, C] \sim [p, B \; ; \; 1-p, C]$

- Monotonicity   $A \succ B \Rightarrow [p > q \Leftrightarrow [p, A \; ; \; 1-p, B] \succ [q, A \; ; \; 1-q, B]]$

- Decomposability   $[p, A \; ; \; 1-p, [q, B \; ; \; 1-q, C]] \sim [p, A \; ; \; (1-p)q, B \; ; \; (1-p)(1-q) C]$

if violated, then agent not rational.

Preferences to Utility

- Existence of Utility Function   (not unique)

If agent's preferences obey axioms of utility, then
$$\exists U \to U(A) > U(B) \text{ iff } A \succ B$$
$$U(A) = U(B) \text{ iff } A \sim B$$

- Expected Utility of a lottery

the utility of a lottery is the sum of the prob.
of each outcome times the utility of that outcome

$$U\left([p_1, s_1 ; \ldots ; p_n, s_n]\right) = \sum_i p_i U(s_i)$$

value function or ordinal utility function : preference
ranking on states

utility : lotteries $\to \mathbb{R}$

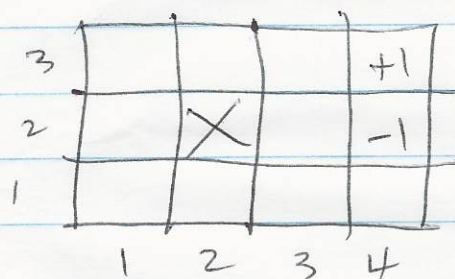preference elicitation : determine agent's utility function

normalized scale   $U(s) = u_\top$  top, best    $= 1$
$U(s) = u_\perp$  bottom, worst $= 0$

assess utility of state $s$ : prob. $p$ $\Rightarrow$ agent is
indifferent in choice between $s$ &
standard lottery $[p, u_\top ; (1-p) u_\perp]$

Chapter 17

sequential decision making in stochastic environment
- utilities
- uncertainty
- sensing
-- search + planning



to make things
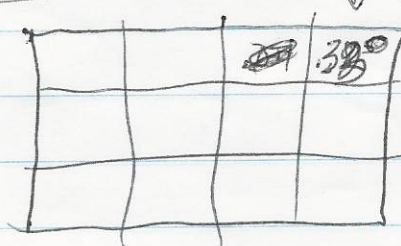easier ↓

(1,1) start state                    } fully observable
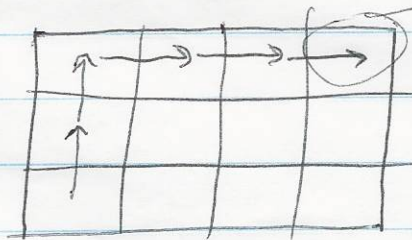(4,2) (4,3) final states             }
Actions(s) = A(s) = {Up, Down, Left, Right}
R((4,2)) = -1    R((4,3)) = 1    R(s) = -0.04    s ≠ (4,2) ∪ (4,3)
Nondeterministic action result
+    .8  prob  direction of selected action
     .1  prob   1 orthog "   "   "   "   "
     .1  prob   other  "   "   "   "   "
                                              ↓ .8 ↑.5 = .3277

Up    .1=Left    .8=Up    .1 = Right

.08
.08
.08

(1,1)

Up    .1=Left    .8=Up    .1=Right    .1=Left    .8=vp    .1=Right    .1=Left    .8=vp    .1=Right

(1,1)    (1,2)    (2,1)    (1,2)    (1,3)    (1,2)    (1,1)    (2,1)    (3,1)

Right    .1=vp    .8    .1=D

(1,2)    (2,1)    (1,1)

| 1,1 | 1,2 | 2,1 | 1,3 | 3,1 |
|-----|-----|-----|-----|-----|
| 2   | .1*.8 | .1*.1 | .8*.8 | .1*.1 |
| .1  | .8*.1 | .1*.8 |       |       |
| 0   | .8*.1 |       |       |       |

| 0.09 | 0.24 | 0.09 | 0.64 | 0.01 | = Σ = 1 |
|------|------|------|------|------|
| 0.02 | 0.24 | 0.09 | 0.64 | 0.01 |

64
24
9
1
2
0

| 0.64 |      |      |
|------|------|------|
| 0.24 |      |      |
| 0.04 | 0.09 | 0.01 |

<u>transition model</u> : outcome of each action in each state

$$P(s' | s, a) \qquad (\text{Markovian : only depends on s})$$

| s' : | 0 | 0 | 0 | 0 |
|------|---|---|---|---|
| (1,2) | .8 | .1 | .1 | 0 |
| (2,1) | .1 | 0 | .8 | .1 |
| (1,1) | .1 | .8 | .1 | .9 |

U    L    R    D
a

## utility function

in each state $s$
agent receives a reward $R(s)$ (pos or neg)

e.g., $-0.04$ in all states except terminal
$$R((4,3)) = +1 \qquad R((4,2)) = -1$$

utility of environment history is sum of rewards

## Markov Decision Process (MDP):

a sequential decision problem
- fully observable
- stochastic environment
- Markovian transition model
- additive rewards

## consists of:

set of states        ($s_0 \equiv$ initial state)
set of actions      ACTIONS($s$)
transition model    $P(s' | s, a)$
reward function     $R(s)$        [could be $R(s, a, s')$]

<u>solution</u> : specify what agent should do in any state
called a <u>policy</u> $\pi$

action from $\pi$ is $\pi(s)$

<u>complete policy</u> : knows what to do in any state

<u>policy quality</u> : measured by <u>expected utility</u>

<u>optimal policy</u> $(\pi^*)$ : highest expected utility

| → | → | → | +1 |
|---|---|---|----|
| ↑ | ✗ | ↑ | −1 |
| ↑ | ↑ | ← | ← |

optimal policy
(for $R(s) = -0.04$

(Note action in (4,1) has 0.1 prob of going to −1
" " " (3,2) " " " " " "

Is that best? Yes, for balance of risk + reward!

if $R(s) \leq -1.628$ then go to nearest exit
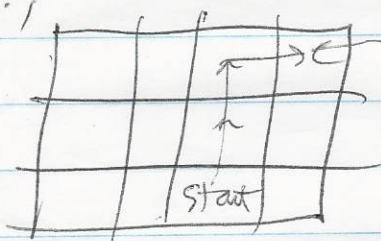
$> 0$ stay on board

<u>Need an algorithm</u>!

Finite horizon or infinite?

game over in fixed time, N

$$U_h([s_0, s_1, ..., s_{N+k}]) = U_h([s_0, s_1, ..., s_N])$$

↖ history

E.g.,



over policy

↓ N=3

N=100 ⇒ go around

nonstationary: policy may change given more or less finite time

stationary: no fixed deadline

Assume: preferences between state sequences are stationary

sequences, $[s_0, s_1, ...]$
$[s_0, s_1', ...]$

if 2 states have same first state, then their
preference is same as $[s_1, s_2, ...]$ & $[s_1', s_2', ...]$

then ∃ 2 ways to assign utilities to sequences:

1. Additive Rewards $U_h([s_0, s_1, ...]) = R(s_0) + R(s_1) + ...$

2. Discounted Rewards: $U_h([s_0, s_1, ...]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + ...$

where $0 < \gamma < 1$ (discount factor)

if can have infinite sequence, then undiscountd
rewards will be infinite.

To avoid:

1. Discounted:

$$U_h([s_0, s_1, s_2, \dots]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \le \sum_{t=0}^{\infty} \gamma^t R_{max} = \frac{R_{max}}{(1-\gamma)}$$

2. If guaranteed to go to terminal state, then ok
3. average reward per time step

We use 1.

Compare policies by comparing expected utilities.

$$U^{\pi}(s) = E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

$$\pi_s^* = \underset{\pi}{\text{argmax}} \ U^{\pi}(s) \quad \left\} \begin{array}{l} \text{is a policy for every state} \\ \text{independent of start state} \end{array} \right.$$

utility of a state is :  $U^{\pi^*}(s) \overset{\text{write as}}{=} U(s)$

| 0.812 | 0.868 | 0.918 | +1 |
| 0.762 | X | 0.660 | -1 |
| 0.705 | 0.655 | 0.611 | 0.388 |

$$\pi^*(s) = \underset{a \in A(s)}{\text{argmax}} \sum_{s'} P(s'|s,a) U(s')$$

# Value Iteration

## Bellman equation for utilities

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) U(s')$$

under-braced: discounted utility of next state

| 9 | 10 | 11 | 12 |
|---|----|----|----|
| 5 | 6 | 7 | 8 |
| 1,1 | 2 | 3 | 4 |

$= 1$

use known utilities to
show $U_p$ is best action

$\downarrow 0.7456$

$U(1,1) = -0.04 + \gamma \max \left[ \overset{0.762}{0.8\, U(1,2)} + 0.1 \overset{0.655}{U(2,1)} + 0.1 \overset{0.705}{U(1,1)} \right.^{\frac{4}{2}}$

$0.705 \qquad 0.762$

$\boxed{0.7107}\; 0.9\, U(1,1) + 0.1\, U(1,2), \qquad\qquad Left$

$\boxed{0.700}\; 0.9\,\overset{0.705}{U(1,1)} + 0.1\,\overset{0.655}{U(2,1)}, \qquad\qquad Down$

$\boxed{0.6707}\; 0.8\,\overset{0.655}{U(2,1)} + 0.1\,\overset{0.762}{U(1,2)} + 0.1\,\overset{0.705}{U(1,1)} \left.\right]\; Right$

## Algorithm

n equations + n unknowns
non linear equations (max)

Use iterative method (Bellman update):

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) U_i(s')$$

function Value-Iteration(mdp, ε) returns a utility function

    inputs: mdp : $S, A, P, R, \gamma$
         ε max error

    local vars  $U, U'$  initially 0
         $\delta$  max change in utility of any state
              in iteration

repeat

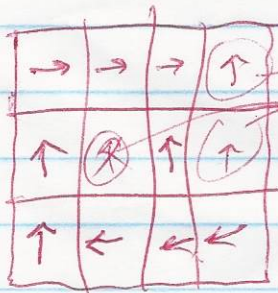    $U \leftarrow U'$; $\delta \leftarrow 0$
    for each state $s$ in $S$ to

        $U'[s] \leftarrow R(s) + \gamma \max\limits_{a \in A(s)} \sum\limits_{s'} P(s' | s, a) U[s']$

        if $|U'[s] - U[s]| > \delta$
        then $\delta \leftarrow |U'[s] - U[s]|$
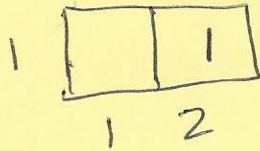until $\delta < \epsilon(1-\gamma)/\gamma$
return $U$


$[S, A, R, P, U, U+] = CS638\_run\_value\_iteration$
           $(0.999999, \overset{3}{50})$;



      initialized to Up & don't change

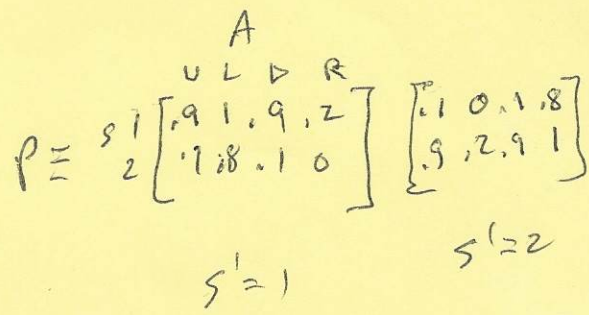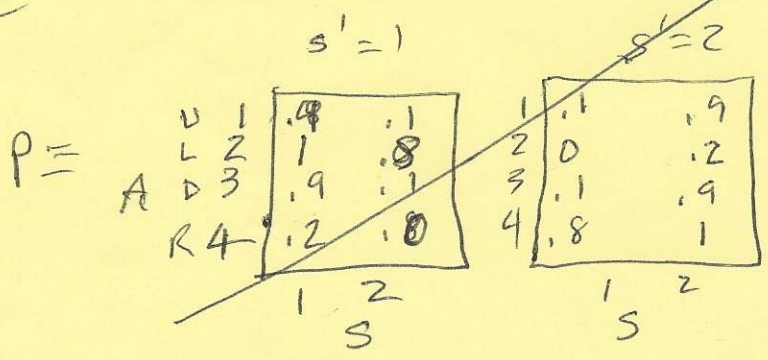          $p = CS6380\_MDP\_policy(S, A, P, U)$

World



$$S = \{[1;1], [2;1]\} \equiv \{1, 2\}$$
$$A = \{Up, Left, Down, Right\} \equiv \{1, 2, 3, 4\} \quad \text{order matters!}$$
$$P = \text{standard} \quad .8, .1, .1, 0$$
$$\gamma = 0.999999$$
$$R = -0.04 \quad \text{at } [1;1]$$
$$\varepsilon = 0.1$$

$$P = $$



$$P \equiv \begin{array}{c} \quad A \\ \begin{array}{c} U\ L\ D\ R \end{array} \\ \begin{array}{c} s\ 1 \\ 2 \end{array} \begin{bmatrix} .9 & 1 & .9 & .2 \\ .7 & .8 & .1 & 0 \end{bmatrix} \end{array} \begin{bmatrix} .1 & 0 & .1 & .8 \\ .9 & .2 & .9 & 1 \end{bmatrix}$$

$$s'=1 \qquad\qquad s'=2$$

$$U \leftarrow 0$$
$$U' \leftarrow 0$$
$$\delta \leftarrow 0$$

repeat
$$U \leftarrow U'$$
$$\delta \leftarrow 0$$
$$\cancel{\times} \quad \frac{\delta \text{ loop}}{s\ a} = 1$$

$$U'(1) = -0.04 + 0.999999 \max_{a \in A(s)} \sum_{s'} P(s' | s_0, a) U(s')$$

* a loop

  a = 1 ~~a, 2~~

  * s' loop

    $s' = 1$

    $V_1 = 0 = P(1|1,1) U(1)$
    $\quad \quad \quad .9 \quad \quad = 0$

    $s' = 2$

    $V_2 = .1 = P(2|1,1) U(2)$
    $\quad \quad \quad .1 \quad \quad = 1$

    end s' loop

    a1. val = .1

  a = 2

  s' loop $s' = 1$

    $V_1 = P(1|1,2) U(1) = 0$
    $\quad \quad = 0$

    $s' = 2$

    $V_2 = P(2|1,2) U(2) = 0$
    $\quad \quad = 0 \quad \quad = 1$

    end s' loop

  a2_val = 0

  a = 3

    s' loop

    $s' = 1$

    $V_1 = P(1|1,3) U(1) = 0$
    $\quad \quad \quad \quad = 0$

    $V_2 = P(2|1,3) U(2) = .1$
    $\quad \quad .1 \quad \quad = 1$

  a3_val = .1

  a = 4

  [ $V_1 = P(1|1,4) U(1) = 0$

  a4. val = .8 $\quad V_2 = P(2|1,4) U(1) = .8$
  $\quad \quad \quad .8 \quad \quad = 1$

$U'(1) \leftarrow .76$

if $|U'(1) - U(1)| > \delta \equiv \quad 0.76 > 0$

$\quad \delta = 0.76$

and

until $\quad \delta < \epsilon (1-\gamma)/\gamma \quad\quad 0.76 < 10^{-7}$

loop again

$U(1) \leftarrow 0.76$

$\delta \leftarrow 0 \quad \cancel{\text{fonds}} \; U'(1) = 0.76$

$\cancel{\text{tends}}$

$s = 1$

$U'(1) \leftarrow 0.04 + 0.999999 \; \underset{a \in A}{max} \; \underset{s'}{\sum} P(s' | s, a) U(s')$

$a = 1$

$\quad s' = 1$

$\quad\quad V_1 = P(1 | 1, 1) * 0.76 \quad\quad 0.684 \;\Big\}\; 0.76$

$\quad\quad\quad .9 * .76$

$\quad s' = 2$

$\quad\quad V_2 = P(2 | 1, 1) * 0.76 \quad\quad = 0.076$

$\quad\quad\quad .1$

$a = 2$

$\quad s' = 1$

$\quad\quad V_1 = P(1 | 1, 2) * .76 \quad = 0.76 \;\Big\}\; 0.76$

$\quad\quad\quad 1 * .76$

$\quad s' = 2$

$\quad\quad V_2 = P(2 | 1, 2) * 0.76 \quad = 0$

$\quad\quad\quad = 0$

$a = 3$

$\quad s' = 1$

$\quad\quad V_1 = P(1 | 1, 3) * .76 \quad 0.684 \;\Big\}\; 0.76$

$\quad\quad\quad .9 * 0.76$

$\quad\quad V_2 = P(2 | 1, 3) * 1 \quad .076$

$a = 4$

$s' = 1$

$v_1 = P(1|1,4) * 0.76$     $0.152$

$\quad .2 * 0.76$          $\Big\}\ .952$

$s' = 2$

$v_2 = P(2|1,4) * \cancel{0.76} = .8$

$\quad .8$

$U'(1) = -0.04 + 0.952 = 0.912$

$\vdots$

— Policy selection given utilities & transition probabilities

$$\Pi^*(s) = \text{argmax} \left\{ \sum_{s'} \left[ P(s'|s,a) U(s') \right] \right\}$$

$a \in A(s)$

pi_star

$$\vdots \begin{bmatrix} \\ \\ \end{bmatrix}$$
$n$

for every state $s$

best_action = 0
best_val = -Inf

for every action $a$

a_sum ← 0

for every state $s'$

add in
← _____ $P(s'|s,a) U(s')$

a_sum = a_sum +

end

pi_star(s) = ?        best action

$$U = \begin{bmatrix} 0.95 \\ 1.0 \end{bmatrix} \qquad P \text{ as before}$$

$S = 1$

best-action = 0

best-val = ~ Inf

a. loop

V

$a = 1$

$a\text{-sum} = 0$

s' loop

$s' = 1$

$a\text{-sum} = P(1|1,1)U(1) = .9 * 0.95 = 0.855$

$s' = 2$

$a\text{-sum} = 0.855 + P(2|1,1)1 = 0.855 + .1 * .95 = .1 = 0.955$

L    $\underline{a = 2}$

$P(1|1,2)U(1) + P(2|1,2) \cdot 1 = 0.095$

   $.1 * 0.95 \quad + \quad 0$

D    $\underline{a = 3}$

$P(1|1,3)U(1) + P(2|1,3) \cdot 1 = 0.955$

R    $\underline{a = 4}$

$P(1|1,4)0.95 + P(2|1,4) \cdot 1 = 0.99$

   $.2 * 0.95 + .8 * 1$

Policy Iteration

given an initial policy $\pi_0$
alternate between:

- Policy evaluation:

   given a policy $\pi_i$, calculate $U_i = U^{\pi_i}$

- Policy improvement

   calculate a new MEU policy $\pi_{i+1}$
   using one step look-ahead based on $U_i$

until no improvement
There are only finitely many policies for finite state space

Policy evaluation   simpler than standard Bellman
                    because no $\underline{max}$ : given policy fixes action

$$U_i(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s')$$

linear set of equations

or iterate $k$ times

$*$ use this

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum P(s' | s, \pi_i(s)) U_i(s')$$

p. 657

function Policy-iteration (mdp) returns a policy
  inputs : mdp: S, A, P
  local var's : U   utilities, initially 0
                π   a policy, initially random


repeat

  U ← Policy-evaluation (π, U, mdp)
  unchanged ← true;
  for each state s in S do
    if $\max\limits_{a \in A(s)} \sum\limits_{s'} P(s'|s,a) U[s'] > \sum\limits_{s'} P(s'|s, \pi[s]) U[s']$

    then do
      $\pi[s] \leftarrow \text{argmax}\limits_{a \in A(s)} \sum\limits_{s'} P(s'|s,a) U[s']$
      unchanged ← false
    endif
  end for
until unchanged == true
return π


p-pi = CS6380_mdp_policy_iteration $(S, A, P, R, k, \gamma)$

R = R+10

| ↓ | ← | ← | 1000 |
| ↓ | ✗ | ← | ~100 |
| ↓ | ← | ← | ↓ |

```
      1
      ↑
2 ←  ↓  → 4
      ↓
      3
```