# Sequential decision making under uncertainty

### Jiří Kléma

Department of Computer Science, Czech Technical University in Prague



http://cw.felk.cvut.cz/doku.php/courses/a4b33zui/start

### **Agenda**

- Previous lecture: individual rational decisions under uncertainty
  - uncertainty = stochastic action outcomes lottery, expected utility,
- how to optimally choose whole sequences of actions?
  - repeated decisions based on uncertain or incomplete information,
  - prizes/rewards typically delayed,
  - world does not have to be fully observable,
- Markov decision process
  - introduces Markov assumption/property process with a limited memory,
  - works with stationarity and observability assumptions too,
  - world/environment well defined by its transition and reward functions,
- generalization
  - POMDP the world is partially observable only,
  - reinforcement learning no environment model available, no state-transition and reward functions.

### Markov process

- random process, prob of visiting future states given by recent states only,
- distant past is irrelevant provided that we know the recent past,
- Markov chain
  - discrete random process with Markov property,
  - chain order m gives how many past states we need to concern

$$P(X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_1 = x_1) = P(X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_{n-m} = x_{n-m})$$

- most often 1st order models,
- commonly together with stationarity assumption (time invariance)

$$Pr(X_{n+1} = x | X_n = y) = Pr(X_n = x | X_{n-1} = y)$$

- examples of Markov chains
  - coin tosses HTHHHT . . .
    - \* degenerate (zero order) Markov chain,
  - weather observed every day at noon SSSCRRCSRR . . .
    - \* categorized (S)unny, (C)loudy, (R)ain, the order is unknown.

### Markov chain – weather example

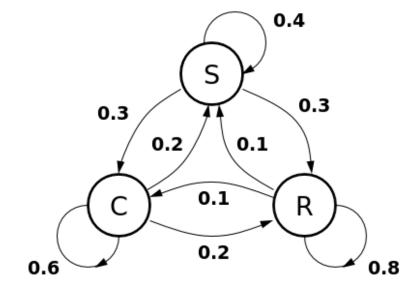
- How to build a model based on the observed sequence?
  - let us have a sequence of length 41, assume 1st order,
     SSCCRRRRRRRSSSCRRRRRRRCCCSCSSRRSRRRCSCCC
  - model = transition matrix.

$$S_{t+1}C_{t+1}R_{t+1}$$

$$S_{t} \begin{bmatrix} 4 & 3 & 3 \\ 2 & 6 & 2 \\ R_{t} \end{bmatrix}$$

$$R_{t} \begin{bmatrix} 2 & 6 & 2 \\ 2 & 2 & 16 \end{bmatrix}$$

$$A = \{a_{ij}\} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$



### Markov chain – weather example

- Which questions can be answered with the Markov model of the sequence?
  - 1. sunny today, probability that weather will be SSRRSCS for next 7 days?
  - 2. model in a known state, prob that the state will not change for d days?
  - 3. what is the expected value d in the individual states?

### Markov chain – weather example

- Which questions can be answered with the Markov model of the sequence?
  - 1. sunny today, probability that weather will be SSRRSCS for next 7 days?
  - 2. model in a known state, prob that the state will not change for d days?
  - 3. what is the expected value d in the individual states?
- Solution

1. 
$$P(O|M) = P(S, S, S, R, R, S, C, S|M) =$$
  
=  $P(S) P(S|S) P(S|S) P(R|S) P(R|R) P(S|R) P(C|S) P(S|C) =$   
=  $1 \times 0.4 \times 0.4 \times 0.3 \times 0.8 \times 0.1 \times 0.3 \times 0.2 = 2.3 \times 10^{-4}$ 

2. 
$$O = \{\underbrace{Q_i, Q_i, \dots, Q_i}_{d}, Q_j \neq Q_i\}, P(O|M, q_1 = Q_i) = a_{ii}^{d-1}(1 - a_{ii}) = p_i(d),$$

3. 
$$\bar{d}_i = \sum_{d=1}^{\infty} dp_i(d) = \sum_{d=1}^{\infty} da_{ii}^{d-1}(1 - a_{ii}) = \frac{1}{1 - a_{ii}}$$
, (sum of arithmetic-geometric series:  $\sum_{k=0}^{\infty} kr^{k-1} = \frac{1}{(1-r)^2}$ ),  $d_S = 1.67, d_C = 2.5, d_R = 5.$ 

### Sequential decision making under uncertainty

- commonly it is more steps/actions needed to reach the goal,
- let us assume
  - non-deterministic environment (actions with uncertain outcomes),
  - the goal state replaced by the aim of maximizing cumulative reward,
- the sequence of actions cannot be found by classical planning
  - rational agent re-examines its steps during the process of solution (execution of actions),
  - next action depends on current observations,
  - current observations depend on current state (= previous actions),
- solution
  - agent evaluates states instead of direct creation of action sequences,
  - in each state we take the action leading to successor states with highest value.

### Basic concepts, problem definition

### lacktriangle Reward $R_t$

— simple sum of immediate rewards obtained per episode:

$$R_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_T$$

- discounted sum for infinite processes ( $\gamma$  is discount rate,  $0 \le \gamma \le 1$ ):  $R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ 

- Policy  $\pi_t(s, a)$ 
  - is a mapping between states and actions,
  - it gives probability that action a will be executed in state s,
  - optimal policy  $\pi^*$  maximizes the total reward  $R_t$ ,

# Basic concepts, problem definition

- State value  $V^{\pi}(s)$ 
  - expected (cumulative) reward for following policy  $\pi$  starting from state s

$$V^{\pi}(s) = E_{\pi}\{R_t \mid s_t = s\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s \right\}$$

- Action value  $Q^{\pi}(s,a)$ 
  - expected (cumulative) reward starting from state s, taking action a and thereafter following  $\pi$

$$Q^{\pi}(s, a) = E_{\pi}\{R_t | s_t = s, a_t = a\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\}$$

• Goal: find  $\pi^*$  (altogether with  $V^*$ ,  $Q^*$  that serve as means).

### Sequential decision making as finite MDP

- Finite Markov Decision Process (MDP)
  - Markov assumption + the sets of states S and actions A are finite,
  - $-MDP = \{S,A,P,R\}$ , can be written as a transition graph,
  - -P transition probability, R reward function,

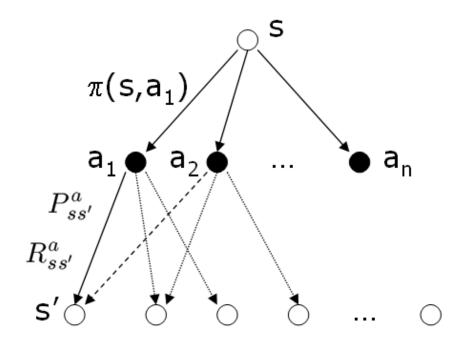
$$P_{ss'}^{a} = Pr\{s_{t+1} = s' \mid s_t = s, a_t = a\}$$

$$R_{ss'}^{a} = E\{r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'\}$$

- this definition leads to particular values of V and Q,
- implicit assumptions
  - environment is observable (the current state is always known),
  - environment is describable (P and R known),
  - counter example: blackjack card game (reaching P and R a part of solution).

# Sequential decision making as finite MDP

- How to obtain state values from known environment and policy?
  - by transition to the recursive V definition,
  - state value = immediate reward for action execution + expected reward for development of possible successor states.



Sutton, Barto: Reinforcement Learning: An Introduction.

# Recursive V definition (Bellman equation)

$$V^{\pi}(s) = E_{\pi} \{ R_{t} \mid s_{t} = s \} =$$

$$= E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = s \right\} =$$

$$= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} \mid s_{t} = s \right\} =$$

$$= \sum_{a} \pi(s, a) \sum_{s'} P_{ss'}^{a} \left[ R_{ss'}^{a} + \gamma E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} \mid s_{t+1} = s' \right\} \right] =$$

$$= \sum_{a} \pi(s, a) \sum_{s'} P_{ss'}^{a} \left[ R_{ss'}^{a} + \gamma V^{\pi}(s') \right]$$

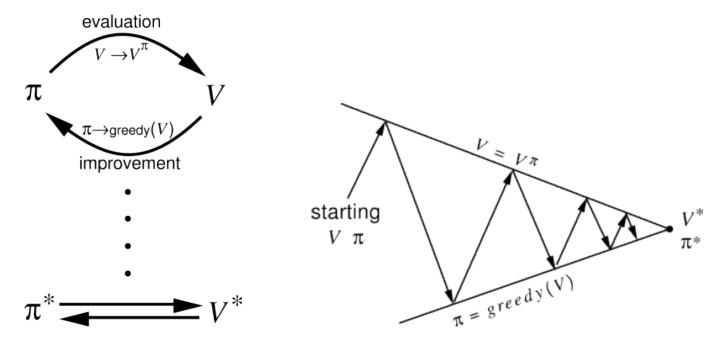
- lacktriangle in the beginning  $V^\pi(s)$ ,  $V^\pi(s')$  and  $\pi(s,a)$  unknown
  - iterative calculation/improvement,
  - bootstrapping effort analogous to one who would lift himself by his own bootstraps.

### **Dynamic programming**

- The basic approach to solve MDP (find  $\pi^*$ )
  - dynamic = iterative procedure
    - \* to find V(s) in step k+1 use V(s') from step k,
  - programming = searching for an acceptable sequence of actions,
- lacktriangle polynomial complexity in the number of states |S| and actions |A|
  - despite the space of policies with cardinality  $|A|^{|S|}$ ,
  - state space search necessarily performs worse,
  - still often intractable for real problems
    - \* unknown process parameters (see reinforcement learning),
    - \* computationally intractable
      - · often too many states,
      - $\cdot$  we cannot iterate systematically asynchronous DP.

### Dependency between policy and value functions

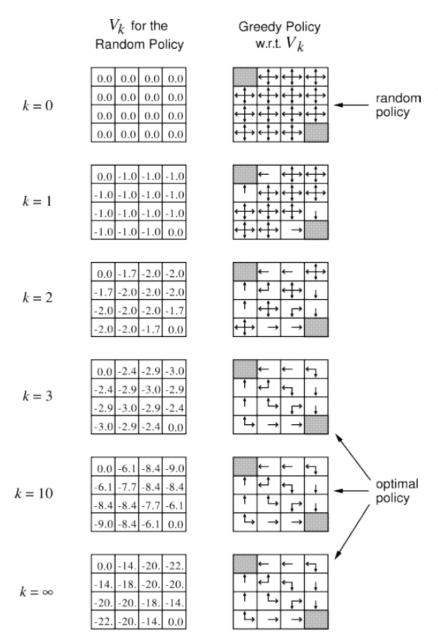
- When solving MDP, one simultaneously and interactively
  - adapts state/action values top the current policy,
  - adapts the policy to maximize reward given the current state/action values.



Sutton, Barto: Reinforcement Learning: An Introduction.

### Policy iteration – PI

- Key idea:  $\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^{\pi^*}$
- 1. policy  $\pi$  evaluation (E step):
  - find state values  $V^{\pi}(s)$ ,
  - start:  $V(s) = 0 \ \forall$  non-terminal states (V known in them),
  - iteration: until V(s) gets steady  $(\max_S |V_{k+1}(s) V_k(s)| < \varepsilon)$ ,
- 2. policy improvement  $\pi \to \pi'$  (I step):
  - adapt to the new state values,
  - deterministic  $\pi$ : in every state takes single action, if  $Q^{\pi}(s,\pi'(s)) \geq V^{\pi}(s)$  for  $\forall s,\,\pi'$  is not worse than  $\pi$ , obviously  $\pi'(s) = \arg\max_a Q^{\pi}(s,a)$ , chooses the currently best action,
  - stochastic  $\pi$ : action selection is driven by a probability distribution, the same logic except for  $Q^{\pi}(s,\pi'(s))=\sum_a \pi'(s,a)Q^{\pi}(s,a)$ ,
- 3. if  $\pi$  and  $\pi'$  differ in at least one state, go to step 1 with  $\pi'$ .



r = -1 on all transitions

### :: Random policy evaluation:

actions

(greedy deterministic to illustrate only)

- $V(s) = \sum_{a} \pi(s, a) \sum_{s'} P^{a}_{ss'} \left[ R^{a}_{ss'} + \gamma V^{\pi}(s') \right]$ ,  $\pi(s, a) = 1/4$ ,  $R^{a}_{ss'} = -1$ ,  $P^{a}_{ss'} = 1$ ,  $\gamma = 1$
- k=0:  $\forall sV(s) = 0$ , no change in terminal states
- k=1:  $V(1) = V(2) = \dots = V(14) = \sum_{a} \pi(s, a) \sum_{s'} P^a_{ss'} \left[ R^a_{ss'} + \gamma V^\pi(s') \right] = 4(1/4 * 1(-1 + 1 * 0)) = -1$
- k=2: V(1) = 1/4(3\*1\*(-1+1(-1))+1\*1(-1+1\*0)) = -7/4 = -1.75
- k=3: V(1) = 1/4(2\*1\*(-1+1(-2)) + 1\*1(-1+1(-1.75)) + 1\*1(-1+1\*0)) = -9.75/4 = -2.44

### Value iteration - VI

- is it necessary to evaluate/know the state values for the given policy perfectly?
  - late iterations often leave policy unchanged,
  - and may spend most of the time of the whole dynamic algorithm,

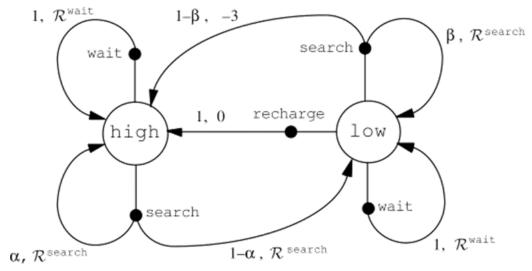
### value iteration

- policy evaluation stopped after first iteration,
- more frequent policy changes,
- in some tasks faster convergence, but does not outperform PI in general,
- in terms of Bellman equation, a new iteration rule originates

$$V^{\pi}(s) = \max_{a} \sum_{s'} P_{ss'}^{a} [R_{ss'}^{a} + \gamma V^{\pi}(s')]$$

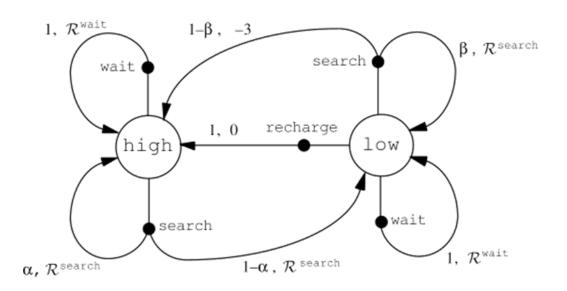
### MDP – recycling robot

- :: Mobile robot that cleans up/collects cans
  - two internal states battery low or high,
  - three actions search for cans, remain stationary, go to home base to recharge,
  - positive reward for each can, negative reward when depleted needing rescue.
- :: Logical goal: collect as many cans as possible without any external aid.
- :: Technical goal: develop a policy that maximizes long term reward.



Sutton, Barto: Reinforcement Learning: An Introduction.

### Recycling robot – DP solution



- Bellman equation:  $V^\pi(s) = \sum_a \pi(s,a) \sum_{s'} P^a_{ss'} \left[ R^a_{ss'} + \gamma V^\pi(s') \right]$
- Iteration equations for particular deterministic action (policy) choices high=h, wait=w, etc.:

$$\begin{split} \pi(h,w) &= 1: & V(h) = Q(h,w) = R^w + \gamma V(h) \\ \pi(h,s) &= 1: & V(h) = Q(h,s) = R^s + \gamma \left[\alpha V(h) + (1-\alpha)V(l)\right] \\ \pi(l,r) &= 1: & V(l) = Q(l,r) = \gamma V(h) \\ \pi(l,w) &= 1: & V(l) = Q(l,w) = R^w + \gamma V(l) \\ \pi(l,s) &= 1: & V(l) = Q(l,s) = \beta R^s - 3(1-\beta) + \gamma \left[\beta V(l) + (1-\beta)V(h)\right] \end{split}$$

# Recycling robot – DP solution

- :: Parameters:  $\alpha=0.95$ ,  $\beta=0.9$ ,  $R^s=2$ ,  $R^w=1$ ,  $\gamma=0.9$ ,  $\varepsilon=0.01$
- :: Method: policy iteration (El cycle)
- 1. Randomly choose a deterministic policy:

$$\pi(low, wait) = \pi(high, wait) = 1$$
,

- 2. set V(low) = V(high) = 0,
- 3. use the iteration equations until V values get steady,
- 4. use evaluations in V to determine optimal actions:

$$V(s) = \max_a Q^{\pi}(s, a), \ \pi'(s) \approx \arg\max_a Q^{\pi}(s, a)$$

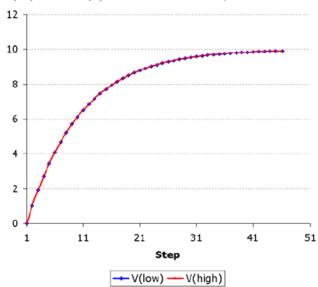
5. in the case of no policy change stop, go to step 2 otherwise.

#### Initialize:

$$\begin{split} \pi(l,w) &= \pi(h,w) = 1 \\ V(h) &= R^w + \gamma V(h), \ V(l) = R^w + \gamma V(l) \end{split}$$

#### Evaluation 1:

$$V(h)=V(l)=10$$
, 46 steps



### Improvement 1:

$$\begin{split} \pi(l,s) &= \pi(h,s) = 1 \\ V(h) &= R^s + \gamma [\alpha V(h) + (1-\alpha)V(l)] \\ V(l) &= \beta R^s - 3(1-\beta) + \gamma [\beta V(l) + (1-\beta)V(h)] \end{split}$$

#### **Evaluation 2**:

$$V(h) = 19, V(l) = 16.8, 52 \text{ steps}$$

### Improvement 2:

$$\pi(l,r) = \pi(h,s) = 1$$

$$V(h) = R^s + \gamma [\alpha V(h) + (1-\alpha)V(l)]$$

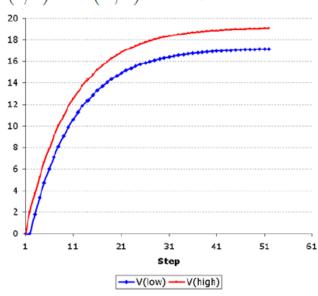
$$V(l) = \gamma V(h)$$

#### **Evaluation 3**:

$$V(h) = 19.1, V(l) = 17.1, 52 \text{ steps}$$

### Improvement 3:

$$\pi(l,r)=\pi(h,s)=1\to \mathsf{STOP}$$



#### SUMMARY:

policy: low  $\rightarrow$  recharge, high  $\rightarrow$  search  $V(h) = 19.1, \ V(l) = 17.1, \ 150$  iterations

# Recycling robot – DP solution

- **::** Parameters:  $\alpha = 0.95$ ,  $\beta = 0.9$ ,  $R^s = 2$ ,  $R^w = 1$ ,  $\gamma = 0.9$ ,  $\varepsilon = 0.01$ ,
- :: Method: value iteration
- 1. Set V(low) = V(high) = 0.
- 2. Use evaluations in V to determine optimal actions:

$$V(s) = max_aQ^{\pi}(s, a), \ \pi'(s) \approx \arg\max_a Q^{\pi}(s, a).$$

- 3. apply once the current best actions and recompute values in V(s),
- 4. in the case of no state value change larger than  $\varepsilon$  stop, go to step 2 otherwise.

http://cw.felk.cvut.cz

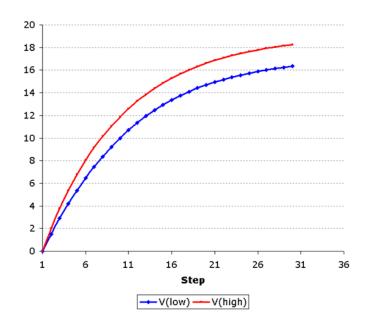
### Recycling robot – DP solution

step 0: V(l) = V(h) = 0,  $\pi(l, s) = \pi(h, s) = 1$ 

step 1: V(l) = 1.5, V(h) = 2,  $\pi(l, s) = \pi(h, s) = 1$ 

step 9: V(l) = 9.2, V(h) = 11.1,  $\pi(l,r) = \pi(h,s) = 1$ , policy change

step 52: V(l)=17.1, V(h)=19.1,  $\pi(l,r)=\pi(h,s)=1$ , all V perturbations smaller than  $\varepsilon$ , STOP



### **SUMMARY:**

policy: low  $\rightarrow$  recharge, high  $\rightarrow$  search V(h)=19.1, V(l)=17.1, 52 iterations

# Why does value iteration certainly converge?

 $lue{}$  contraction c(x)

$$-\exists k \ \forall x_1 x_2 : d(c(x_1) - c(x_2)) \le k d(x_1 - x_2),$$

- -d is a metric (distance function), constant  $0 \le k < 1$ ,
- fixed point  $b_c$ :  $c(b_c) = b_c$ ,  $c(c(\ldots c(x))) = b_c$ ,
- each contraction has only one fixed point,
- example:  $c(x) = \frac{x}{2}$ , d(x, y) = |x y|,  $b_c = 0$ ,
- value iteration equation

$$-V_{i+1}(s) = \max_{a} \sum_{s'} P_{ss'}^{a} [R_{ss'}^{a} + \gamma V_{i}(s')]$$

- can be simplified as  $V_{i+1} \leftarrow BV_i$ ,
- as d we employ max norm  $||V|| = max_s|V(s)|$ ,
- the above defined B is wrt || || contraction (without proof)
  - $-||BV_i BV_i'|| \le \gamma ||V_i V_i'||.$

■ □ □ □ □ □ □ □ □ □ □ □ http://cw.felk.cvut.cz

# Why does value iteration certainly converge?

- lacktriangle provided that B is a contraction wrt  $||\ ||$ 
  - for any pair of state utility vectors it holds

$$||BV_i - BV_i'|| \le \gamma ||V_i - V_i'|| \Rightarrow ||V_{i+1} - V_i|| \le \gamma ||V_i - V_{i-1}||,$$

- \* value iteration converges for  $\gamma < 1$ ,
- the fixed point is the actual state utility vector  $V^st$ 
  - $* ||BV_i V^*|| \le \gamma ||V_i V^*||,$
  - \* converges exponentially with  $\gamma$ .

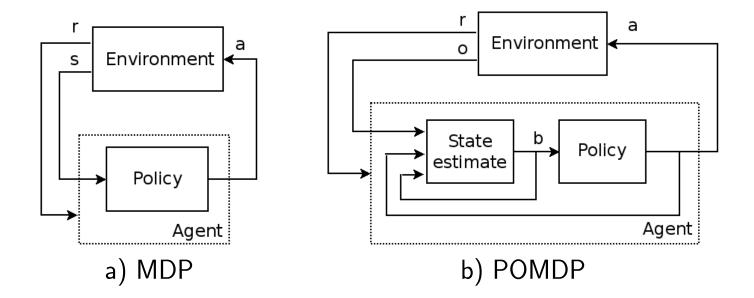
### Partial observability

- MDPs work with the assumption of complete observability
  - assumption that the actual state s is always known is strong and often non realistic,
- partially observable Markov decision process (POMDP)
  - MDP generalization, states can be guessed from observations coupled with them,
  - $POMDP = \{S, A, P, R, O, \Omega\},\$ 
    - \* O is a set of observations,
    - \*  $\Omega$  is a sensoric model defines conditional observation probs,  $\Omega^a_{s'o} = Pr\{o_{t+1} = o \mid s_{t+1} = s', a_t = a\}$
  - instead of s agent internally keeps prob distribution b (belief) across states
    - \* we perform a in unknown s (knowing b(s) only) and observe o then  $b'(s') = \eta \ \Omega^a_{s'o} \sum_{s \in S} P^a_{ss'} b(s)$ ,
    - \*  $\eta$  is a normalization constant such that  $\sum_{s' \in S} b'(s') = 1$ .

■ □ □ □ □ □ □ □ □ □ lattp://cw.felk.cvut.cz

### Partial observability

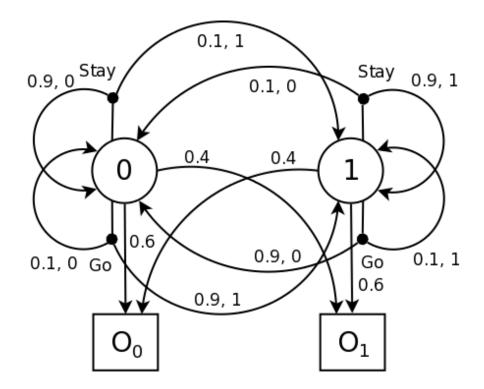
- consequences of partial observability
  - it makes no sense to concern policy  $\pi: S \to A$ , shift to  $\pi: B \to A$ ,
  - commonly computationally intractable, approximate solutions only
    - \* for n states, b is an n-dimensional real vector,
    - \* PSPACE-hard, worse than NP.



# Partial observability – example

:: 
$$S = \{0, 1\}$$
,  $A = \{Stay, Go\}$ ,  $O = \{o_0, o_1\}$ ,  $P(s_{t+1} = x | s_t = x, Stay) = .9$ ,  $P(s_{t+1} = x | s_t = x, Go) = .1$ ,  $Pr(o_0|0) = .6$ ,  $Pr(o_1|1) = .6$ ,  $R(0) = 0$ ,  $R(1) = 1$ ,  $\gamma = 1$ ,

:: Goal: determine  $V^*(b)$  (the main step for finding  $a = \pi^*(b)$ )



■ □ □ □ □ □ □ □ http://cw.felk.cvut.cz

# Partial observability - example

- lacksquare b space is  $1{
  m D} o V(b)$  is a real function of one variable,
- assumed that in near points of b space will be
  - very similar utility and identical policy,
- policy is equivalent to a conditional plan dependent on future observations
  - example, plan of length 2: [Stay], if  $O = o_0$  then Go else Stay,
- let  $\alpha_p(s)$  be the utility of plan p starting from state s
  - then the same plan executed from b has the utility

$$\sum_{s} b(s)\alpha_p(s) = b \cdot \alpha_p$$

- $-\alpha_p$  is a linear function of b (hyperplane for complex spaces),
- optimal policy follows the plan with highest expected utility

$$V(b) = V^{\pi^*}(b) = \max_{p} b \cdot \alpha_p$$

-V(b) is a partially linear function of b.

□ □ □ □ □ □ http://cw.felk.cvut.cz

### Partial observability - example

:: there are two plans of length 1, for them it holds

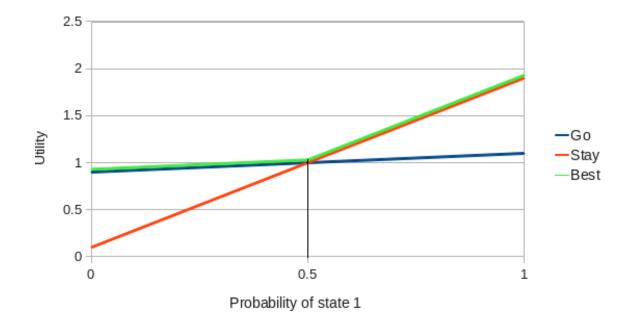
$$\alpha_{[Stay]}(0) = R(0) + \gamma(.9R(0) + .1R(1)) = 0.1$$

$$\alpha_{[Stay]}(1) = R(1) + \gamma(.9R(1) + .1R(0)) = 1.9$$

$$\alpha_{[Go]}(0) = R(0) + \gamma(.9R(1) + .1R(0)) = 0.9$$

$$\alpha_{[Go]}(1) = R(0) + \gamma(.9R(0) + .1R(1)) = 1.1$$

$$\alpha_{[Stay]}(b(1) = 0.3) = .7\alpha_{[Stay]}(0) + .3\alpha_{[Stay]}(1) = 0.64$$



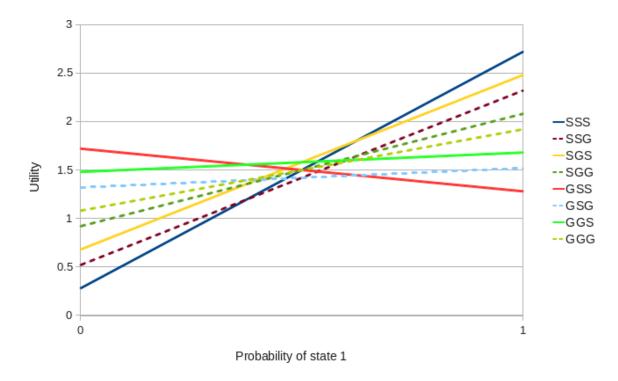
■ □ □ □ □ http://cw.felk.cvut.cz

### Partial observability – example

:: there are 8 plans of length 2 (4 dotted plans strictly dominated by other plans)

 $[Stay, if O = o_0 then Go else Stay]$  encoded as [SGS]

$$\begin{array}{lll} \alpha_{[SSS]}(0) &= R(0) + \gamma(.9\alpha_{[S]}(0) + .1\alpha_{[S]}(1)) &= 0.28 \\ \alpha_{[SSS]}(1) &= R(1) + \gamma(.9\alpha_{[S]}(1) + .1\alpha_{[S]}(0)) &= 2.72 \\ \alpha_{[SGS]}(0) &= R(0) + \gamma(.9(.6\alpha_{[G]}(0) + .4\alpha_{[S]}(0)) + .1(0.4\alpha_{[G]}(1) + .6\alpha_{[S]}(1)) &= 0.68 \\ \alpha_{[SGS]}(1) &= R(1) + \gamma(.9(.4\alpha_{[G]}(1) + .6\alpha_{[S]}(1)) + .1(0.6\alpha_{[G]}(0) + .4\alpha_{[S]}(0)) &= 2.48 \end{array}$$



■ □ □ □ http://cw.felk.cvut.cz

# Partial observability – plans of length d

lacktriangle generalized formula to evaluate plans of length d

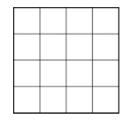
$$\alpha_p(s) = R(s) + \gamma \left( \sum_{s'} P_{ss'}^a \sum_o \Omega_{s'o}^a \alpha_{p.o}(s') \right)$$

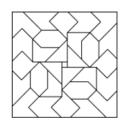
- recursive formula,
- plan p with length d, p.o is its subplan with length d-1 without observation o,
- pruning of dominated plans helps,
- still, worst-case time complexity  $\mathcal{O}(|A|^{|O|^{d-1}})$ .

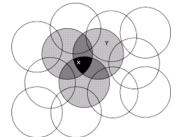
□ □ □ http://cw.felk.cvut.cz

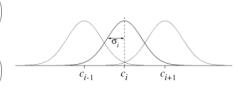
### Generalization and approximation

- up to know trivial demonstrations with limited state and action sets,
- for large state spaces, necessary approximation of value functions
  - learning from examples can be employed,
  - generalization assumes continuous and "reasonable" value functions,
  - states characterized by a parameter/attribute vector  $\phi_s$ ,









- linear approximation with parameters heta(t):  $V_t(s) = heta_t^T \phi_s = \sum_i heta_t(i) \phi_s(i)$ ,
- non linear optimization by a neural network,
- error function to be minimized:  $MSE(\theta_t) = \sum_s P_s \left[V_t^{\pi}(s) V_t(s)\right]^2$ ,  $P_s$  distribution of state weights,  $V_t^{\pi}(s)$  real value,  $V_t(s)$  its approximation,
- regression, gradient optimization, back propagation etc.

### **Summary**

- MDPs allow to search stochastic state spaces
  - computational complexity is increased due to stochasticity,
- problem solving = policy finding
  - policy assigns each state the optimal action, can be stochastic too,
  - basic approaches are policy iteration and value iteration,
  - other choices can be modified iteration approaches, possibly asynchronous,
- techniques similar to MDP
  - POMDP for partially observable environments,
  - RL for environments with unknown models,
- applications
  - agent technology in general, robot control and path planning in robotics,
  - network optimization in telecommunication, game playing.

■ □ http://cw.felk.cvut.cz

### Recommended reading, lecture resources

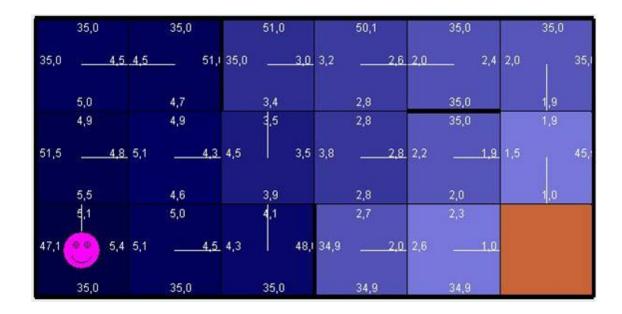
### :: Reading

- Russell, Norvig: Al: A Modern Approach, Making Complex Decisions
  - chapter 17,
  - online on Google books:
     http://books.google.com/books?id=8jZBksh-bUMC,
- Sutton, Barto: Reinforcement Learning: An Introduction
  - MIT Press, Cambridge, 1998,
  - http://www.cs.ualberta.ca/~sutton/book/the-book.html.

http://cw.felk.cvut.cz

### **Demo**

- RL simulator
  - find the optimal path in a maze
  - implemented in Java
  - http://www.cs.cmu.edu/~awm/rlsim/



©Kelkar, Mehta: Robotics Institute, Carnegie Mellon University

Http://cw.felk.cvut.cz