

Using Markov Decision Processes for Policy Enforcement

September 15, 2017

Automated enforcement of privacy policies

- Healthcare industry has had numerous high profile privacy incidents
- Government is increasingly cracking down on such breaches
- Privacy is an *ongoing process*
- Problem: Manual audit log reviews are tedious and expensive

Verifying an organization obeys privacy policies

- Privacy policies restrict the purpose for which some protected information can be used for
- Types of restrictions
 - Prohibitive rules
 - An organization does *not* use certain information *for* a purpose
 - Example: Yahoo!'s privacy policy
 - Exclusivity rules
 - An organization users certain information *only for* a given list of purposes
 - Example: HIPAA

Enforcing privacy policies

- Problem:
 - Tools for enforcement assume that actions they audit are **accurately** labeled with the purpose
 - Tools do not analyze the meaning of *purpose*
- Solution:
 - Need a way to determine if an action could be for a purpose

Purpose in policy enforcement

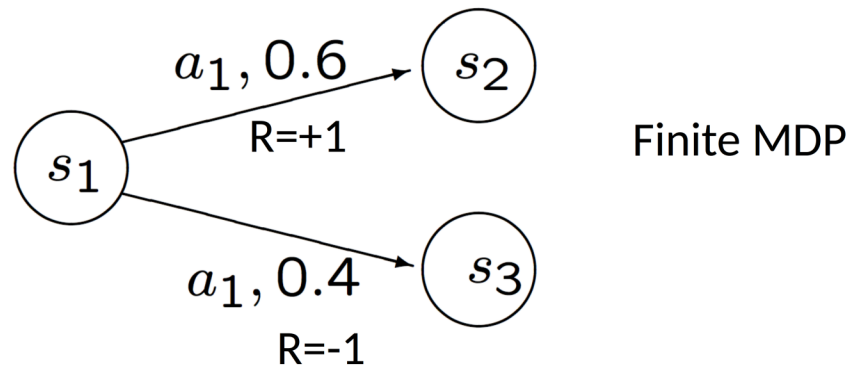
- Planning is central to purpose
 - Cognitive psychology: Purpose is the central determinant of behavior
- Formalizing the relationship between planning and purpose
 - Use MDP to construct all possible behaviors an exclusivity privacy policy allows
 - Auditing an organization:
 - MDP model of allowed behaviors VS
 - Behaviors recorded in the given audit log

Planning for a purpose

- Use an MDP to model the environment the organization being audited operates in
- Why MDP?
 - Good for planning with probabilistic systems
 - Reward function quantifies the degree of progress towards a purpose that taking an action from a state results in
- If the organization being audited logs actions that only help to achieve a particular purpose:
 - Actions correspond to an optimal plan for that MDP
 - Actions are for that purpose

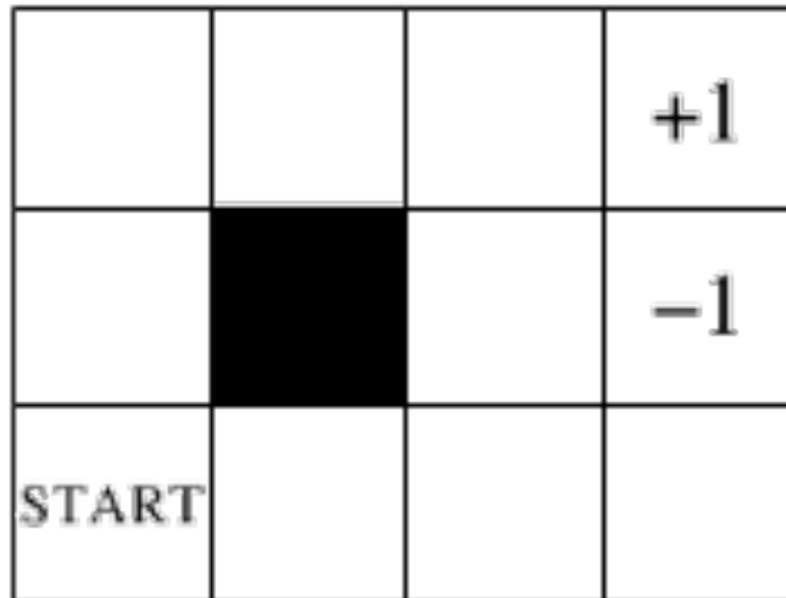
MDP: Formally

- MDP is a tuple $\langle S, A, T, R, \gamma \rangle$
 - S is a finite state space
 - A is a finite action set
 - $T: S \times A \rightarrow D(S)$ is a transition function that given a state and an action, maps to a distribution over states $D(S)$
 - $R: S \times A \rightarrow \mathbb{R}$ is a reward function
 - γ : a discount factor such that $0 < \gamma < 1$



Example process modeled by MDP

- Actions: Left, Right, Up, Down
- Entering top-right corner gives reward of +1 and then takes you to a random state.
- Finding good strategy that maximizes reward.



Goal of an agent using MDP to plan its actions

- $MDP = \langle S, A, T, R, \gamma \rangle$
- ***Maximize its expected total discounted reward***
 - i is a natural number ranging over time modeled as discrete steps
 - R_i is the reward at time i
 - Expectation taken over probabilistic transitions
 - Discount factor γ accounts for preference of people to receive rewards sooner, not later

What is a good strategy?

- What we want a strategy to optimize:
 - Expected reward / time step
 - Expected reward in the first t steps
 - *Maximizing the expected discounted reward*

Maximizing the discounted reward

- $Q(s, a) = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots$

$$V(s) = \max_a Q(s, a)$$

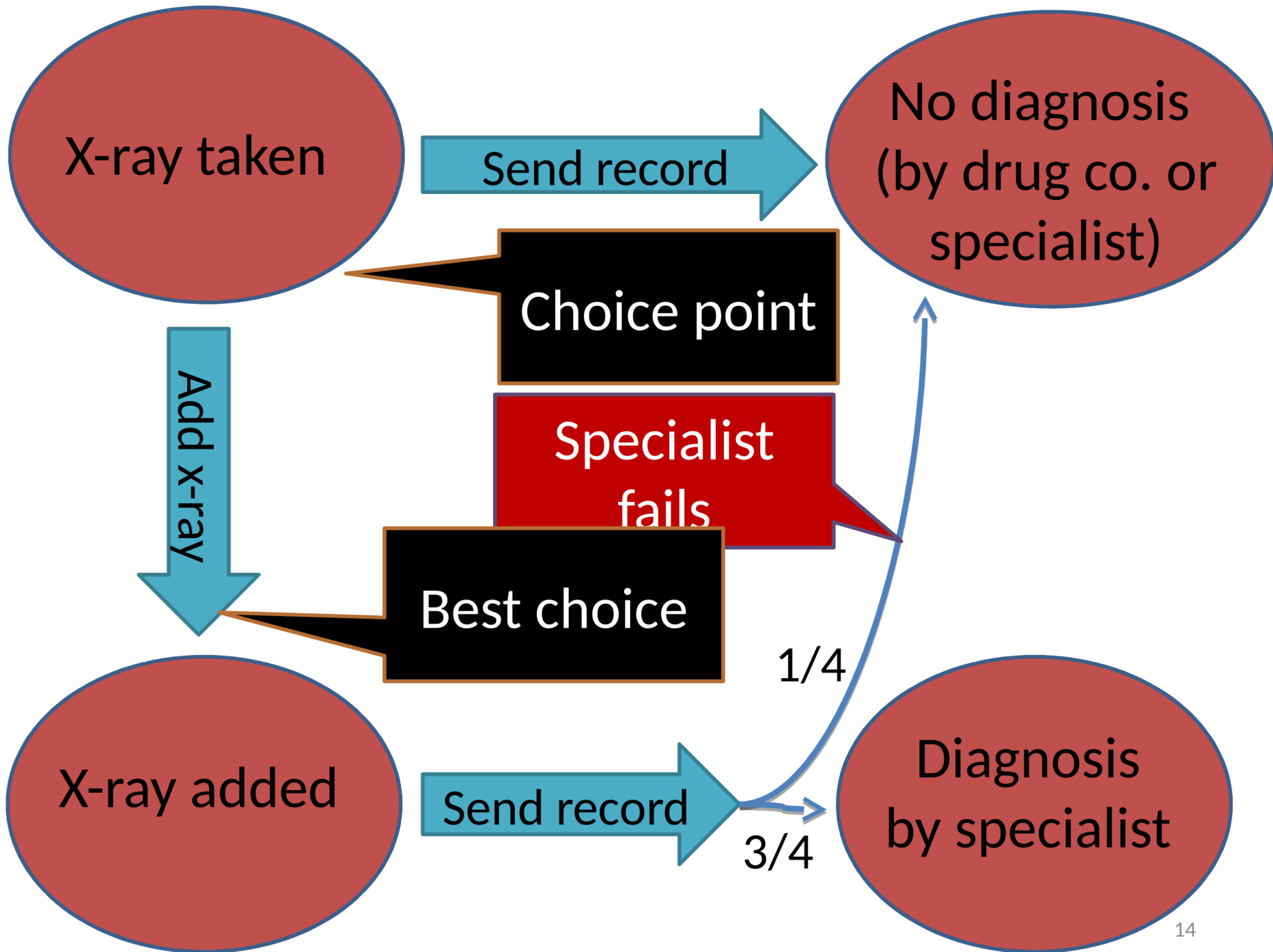
$$V(s) = \max_a [R(s, a) + \gamma \sum_{s'} P(s') V(s')]$$

Solving for the maximum discounted award

- Given:
 - Transition function
 - Reward function
- Possible Solutions
 - Dynamic programming
 - Start with guesses $V_0(s)$ for all states s
 - Update using
$$V_i(s) = \max_a \left[\mathbf{E}[R(s, a)] + \gamma \sum_{s'} \text{Pr}(s') V_{i-1}(s') \right]$$
 - Linear programming
 - Replace “max” with “min” and minimize subject to the constraints

Infinite MDP

There exists an optimal stationary policy, with certain conditions on rewards.



Was that an infinite MDP?

- How can you make it one?
 - Add a stop action that loops at the end

Audit Algorithm

```
AUDIT( $m = \langle \mathcal{S}, \mathcal{A}, t, r, \gamma \rangle$ ,  $b = [s_1, a_1, \dots, s_n, a_n]$ ):  
01  if (IMPOSSIBLE( $m, b$ ))  
02      return true    // behavior impossible for NMDP  
03   $V_m^* := \text{SOLVEMDP}(m)$   
04  for ( $i := 1; i \leq n; i++$ ):  
05      if ( $Q^*(V_m^*, s_i, a_i) < V_m^*(s_i)$ ):  
06          return true    // action suboptimal  
07      if ( $Q^*(V_m^*, s_i, a_i) \leq 0$  and  $a_i \neq \text{Stop}$ ):  
08          return true    // action redundant  
09  return false
```

Figure 2. The algorithm AUDIT

Questions on the homework?