MARKOV DECISION PROCESSES

10-601 Machine Learning Fall 2012 Recitation November 27, 2012 Selen Uguroglu

OUTLINE

- I. MDP overview
- 2. Value Iteration
- 3. Policy Iteration
- 4. Previous Exam Questions

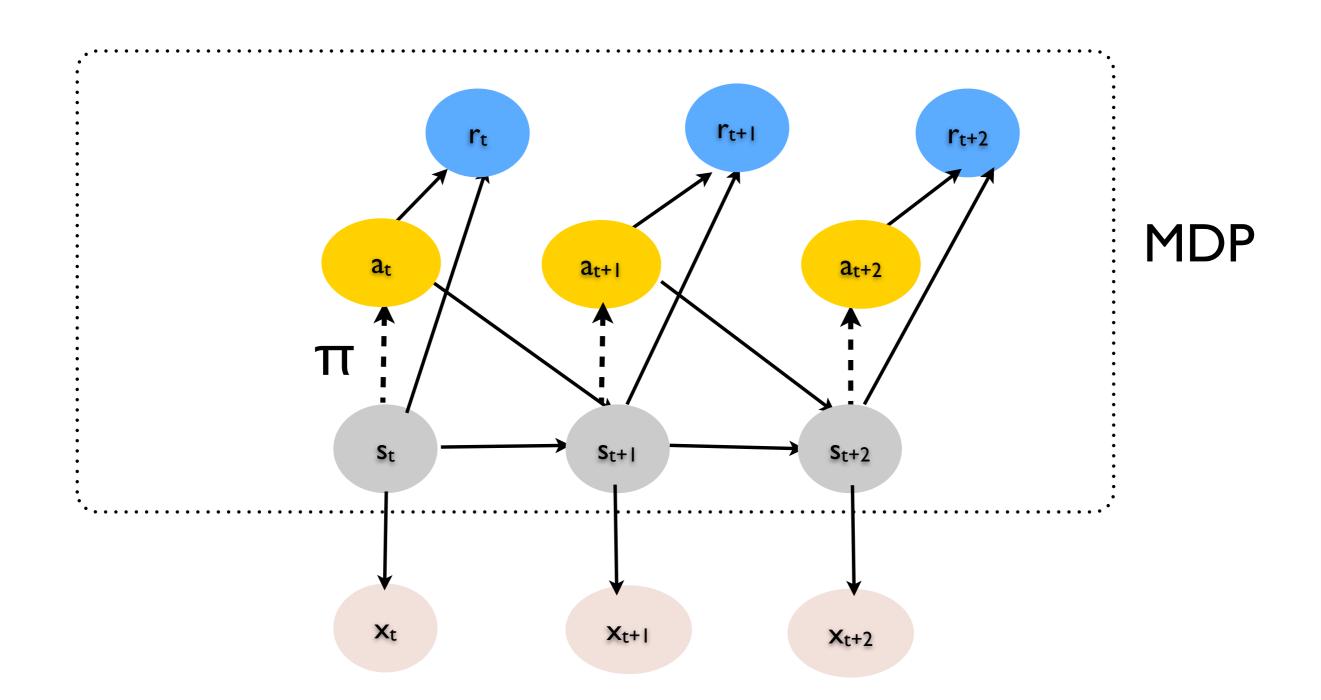
OUTLINE

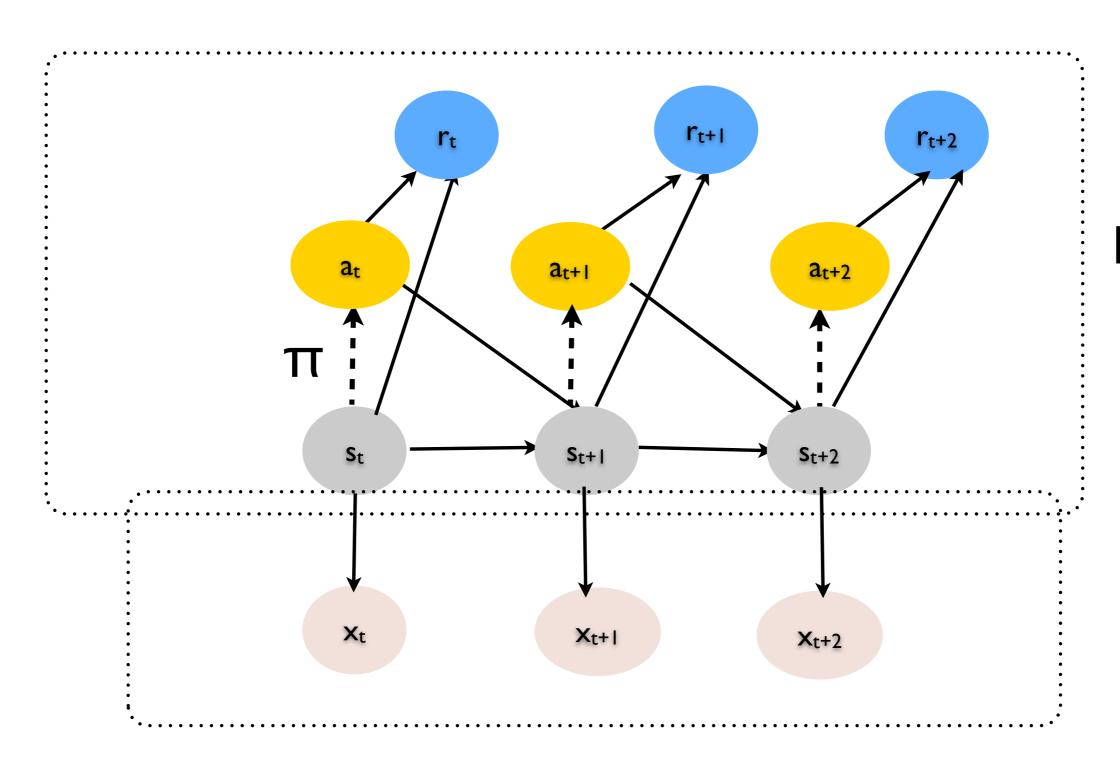
- I. MDP overview
- 2. Value Iteration
- 3. Policy Iteration
- 4. Previous Exam Questions

MARKOV DECISION PROCESSES

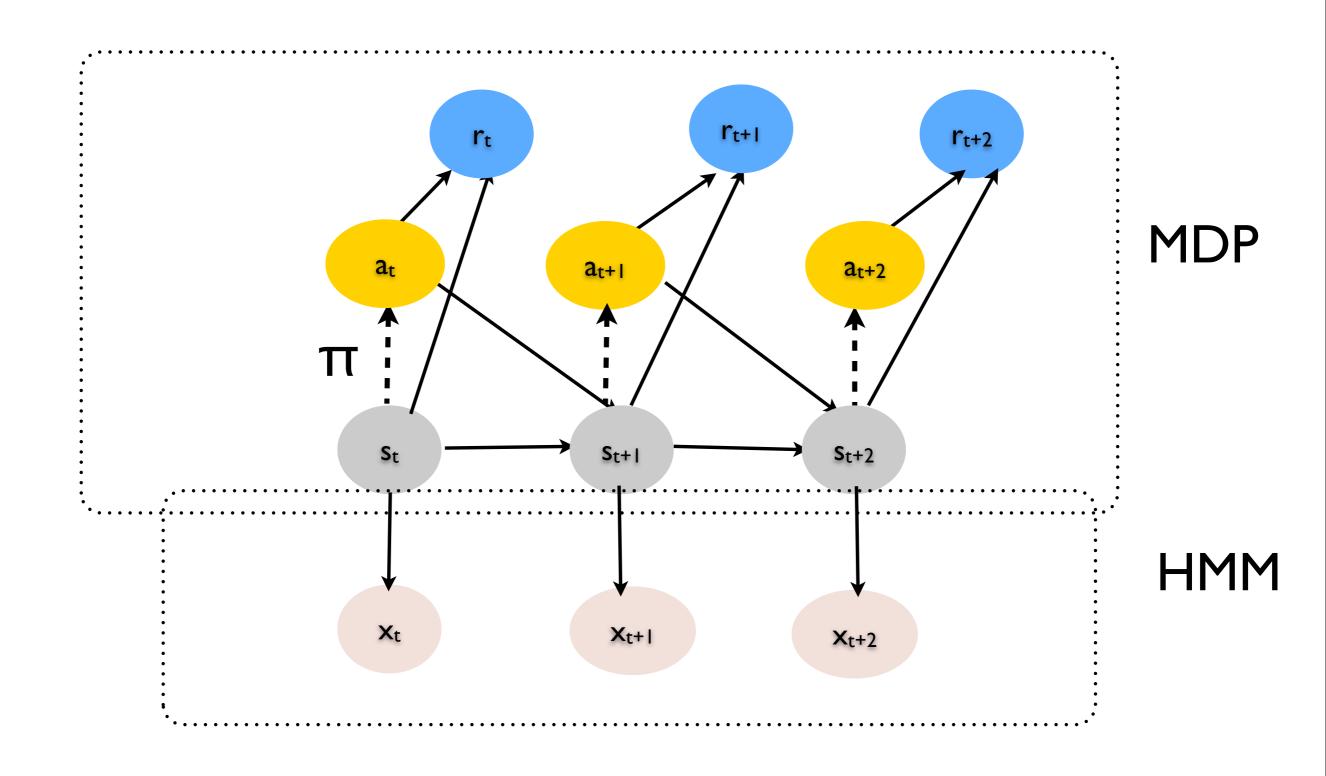
- ▶ MDP is a stochastic process that models the environment under different actions.
- It is defined on states, actions and rewards.
- Agent decides on the action, and receives a reward which depends on the action and states

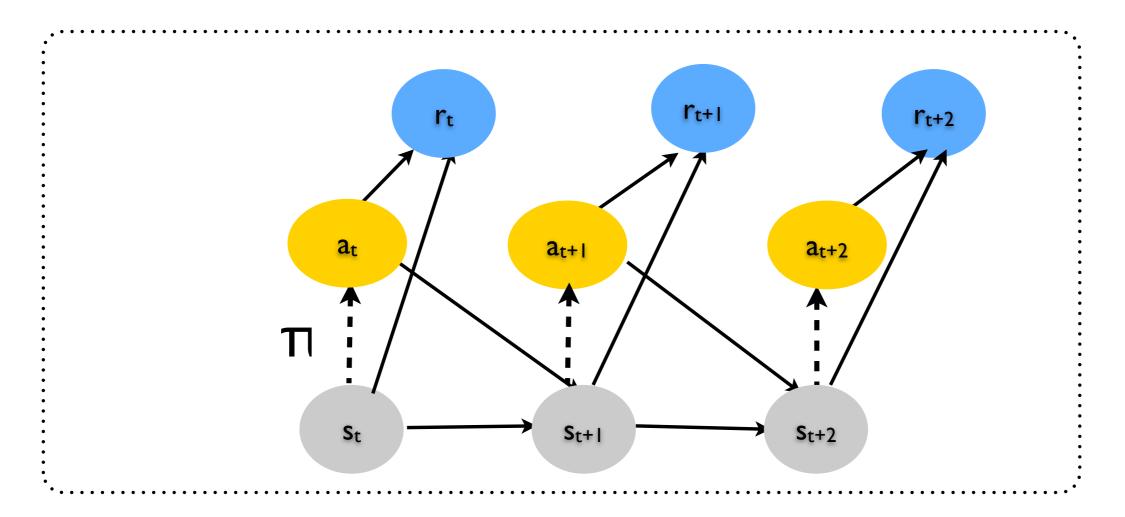
Goal: what action should I take to maximize my chance to win a game?





MDP





MDP

Transition probability

 $P(s_{t+1}|a_{t,s_t})$

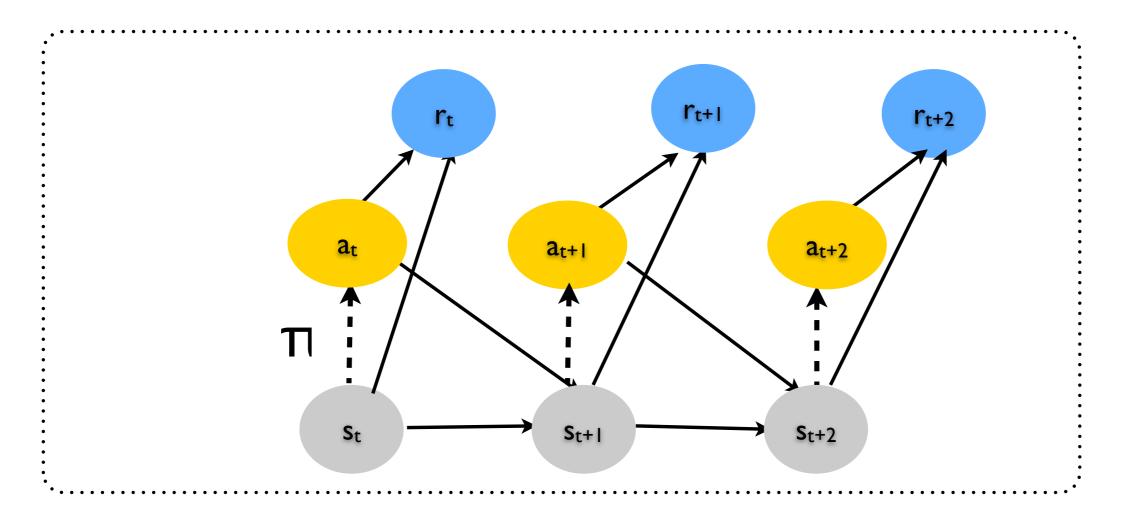
Reward probability

 $P(r_t|a_t, s_t)$

Policy

 $P(a_t|s_t)$

 $\pi(a_t|s_t)$



MDP

Transition probability

 $P(s_{t+1}|a_{t,s_t})$

Pst+1st

Reward probability

 $P(r_t|a_t, s_t)$

Policy

 $P(a_t|s_t)$

 $\pi(a_t|s_t)$

MARKOV DECISION PROCESSES

- Set of states $S = \{s_1, s_2, ..., s_n\}$
- Set of actions $A = \{a_1, a_2, ..., a_n\}$
- Set of rewards $R = \{r_1, r_2, ..., r_n\}$
- Policy π gives an action for each state, π : $S \rightarrow A$ What are the Markov assumptions?

MARKOV DECISION PROCESSES

- Set of states $S = \{s_1, s_2, ..., s_n\}$
- Set of actions $A = \{a_1, a_2, ..., a_n\}$
- Set of rewards $R = \{r_1, r_2, ..., r_n\}$
- Policy π gives an action for each state, π : $S \rightarrow A$ What are the Markov assumptions?

$$P(r_{t}|s_{t},a_{t},s_{t-1},a_{t-1},..) = P(r_{t}|s_{t},a_{t})$$

$$P(s_{t+1}|s_{t},a_{t},s_{t-1},a_{t-1},..) = P(s_{t+1}|s_{t},a_{t})$$

Value function, $V^{\pi}(s)$ is a measure for the expected discounted return:

$$V^{\pi}(s) = E\{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s; \pi\}$$

$$= E\{\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s; \pi\}$$

$$= \sum_{t=0}^{\infty} \gamma^t E\{r_t | s_0 = s; \pi\}$$

Value function, $V^{\pi}(s)$ is a measure for the expected discounted return:

$$V^{\pi}(s) = E\{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s; \pi\}$$

$$= E\{\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s; \pi\}$$

$$= \sum_{t=0}^{\infty} \gamma^t E\{r_t | s_0 = s; \pi\}$$

Discounted by γ exponentially to the future

Value function, $V^{\pi}(s)$ is a measure for the expected discounted return:

$$V^{\pi}(s) = E\{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s; \pi\}$$

$$= E\{\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s; \pi\}$$

$$= \sum_{t=0}^{\infty} \gamma^t E\{r_t | s_0 = s; \pi\}$$

Discounted by γ exponentially to the future

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

Value function, $V^{\pi}(s)$ is a measure for the expected discounted return:

$$V^{\pi}(s) = E\{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s; \pi\}$$

$$= E\{\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s; \pi\}$$

$$= \sum_{t=0}^{\infty} \gamma^t E\{r_t | s_0 = s; \pi\}$$

Discounted by γ exponentially to the future

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

Does it converge?

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

But this for only one policy, π How about the optimal value function, $V^*(s)$?

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

But this for only one policy, π How about the optimal value function, $V^*(s)$? Optimal value function

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

How about the optimal value function, $V^*(s)$?

Optimal value function

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

How about the optimal value function, $V^*(s)$?

Optimal value function

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Optimal policy

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

How about the optimal value function, $V^*(s)$?

Optimal value function

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Optimal policy

$$\pi^*$$
 optimal $\Leftrightarrow \forall_s : V^{\pi^*}(s) = V^*(s)$

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

How about the optimal value function, $V^*(s)$?

Optimal value function

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Optimal policy

$$\pi^*$$
 optimal $\Leftrightarrow \forall_s : V^{\pi^*}(s) = V^*(s)$

Is there really an optimal policy for every MDP?

We can use the recursive property to compute the value function for each state

$$V^{\pi}(s) = E\{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots | s_0 = s; \pi\}$$

$$= E\{r_0 | s_0 = s; \pi\} + \gamma E\{r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | s_0 = s; \pi\}$$

$$= R(\pi(s), s) + \gamma \sum_{s'} P(s' | \pi(s), s) E\{r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | s_1 = s'; \pi\}$$

$$= R(\pi(s), s) + \gamma \sum_{s'} P(s' | \pi(s), s) V^{\pi}(s')$$

For all possible transitions from state s

Recursive property holds for optimal value function V*



Bellman's Equation

With actions

$$V^*(s) = \max_{a} [R(a, s) + \gamma \sum_{s'} P(s' | \pi(s), s) V^*(s')]$$

$$\pi^*(s) = \operatorname{argmax}_a[R(a, s) + \gamma \sum_{s'} P(s' | \pi(s), s) V^*(s')]$$

No actions

$$V^{*}(s) = r_{s} + \gamma \sum_{s'} P(s'|\pi(s), s) V^{*}(s')$$

$$\pi^*(s) = r_s + \gamma \sum_{s'} P(s'|\pi(s), s) V^*(s')$$

Suppose we have n states:

$$V^*(s_1) = r_{s_1} + \gamma \left(p_{s_1 s_1} V^*(s_1) + p_{s_1 s_2} V^*(s_2) + \dots + p_{s_1 s_n} V^*(s_n) \right)$$

$$V^*(s_2) = r_{s_1} + \gamma \left(p_{s_2 s_1} V^*(s_1) + p_{s_1 s_2} V^*(s_2) + \dots + p_{s_n s_n} V^*(s_n) \right)$$

$$V^*(s_n) = r_{s_1} + \gamma \left(p_{s_n s_1} V^*(s_1) + p_{s_n s_2} V^*(s_2) + \dots + p_{s_n s_n} V^*(s_n) \right)$$

It can be solved in n equations in closed form.

But we may not be able to do this every time, so we consult to value/ policy iteration to find the optimal values.

OUTLINE

- I. MDP overview
- 2. Value Iteration
- 3. Policy Iteration
- 4. Previous Exam Questions

VALUE ITERATION

Iterative algorithm, start with $V^0(s_i)$ This could be initialized to 0

$$V^{1}(s_{1}) = r_{s_{1}}$$

$$V^{2}(s_{1}) = r_{s_{1}} + \gamma(\sum_{k} p_{s_{1}s_{k}} V^{1}(s_{k}))$$

$$V^{t+1}(s_{1}) = r_{s_{1}} + \gamma(\sum_{k} p_{s_{1}s_{k}} V^{t}(s_{k}))$$

$$|V_{s_i}^{t+1} - V_{s_i}^t|_{t \to \infty} < \epsilon$$

OUTLINE

- I. MDP overview
- 2. Value Iteration
- 3. Policy Iteration
- 4. Previous Exam Questions

POLICY ITERATION

- Initialization:
 - Randomly choose π_0 set t = 0
- ▶ Policy evaluation: For each s_i, compute V*(s_i)
- Policy update:

$$\pi_t(s_i) = \max_a r_i + \gamma(\sum_j P(s_j|a, s_i)V^*(s_j)$$

▶ If not converged, t = t+I

POLICY ITERATION VS VALUE ITERATION

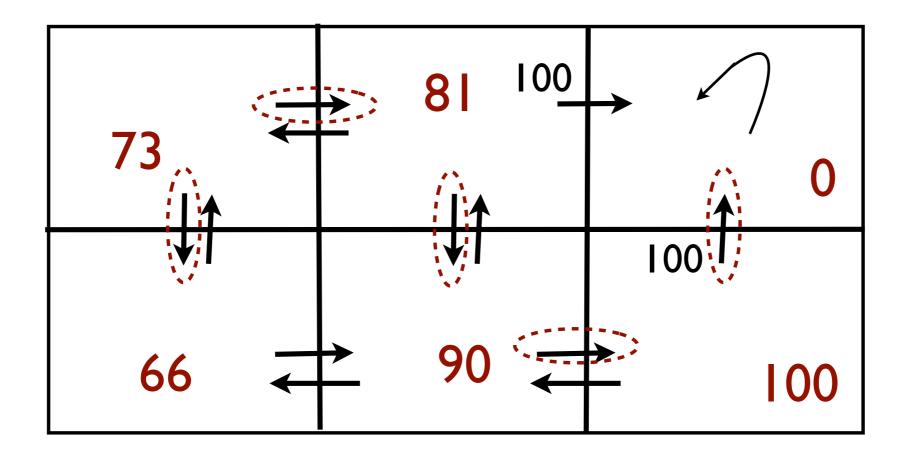
- ▶ Policy iteration is good when the initial policy guess is right
- If we have too many actions value iteration may be slow
- Otherwise value iteration is safer

MARKOV DECISION PROCESSES

- I. MDP overview
- 2. Value Iteration
- 3. Policy Iteration
- 4. Previous Exam Questions

Warm-up

Suppose policy π is circled. Suppose γ is 0.9. What are the V^{π} (s) values? Immediate rewards are written next to transitions, transitions with no immediate reward has 0 value.

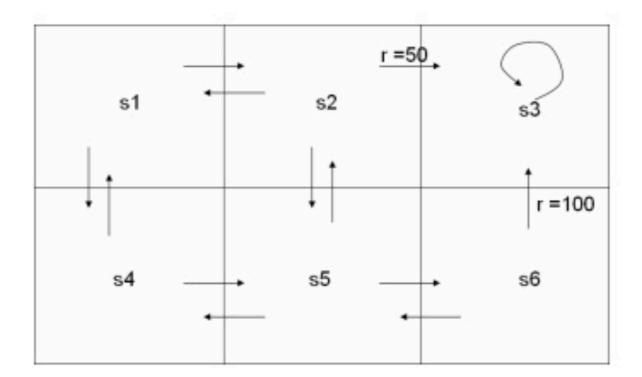


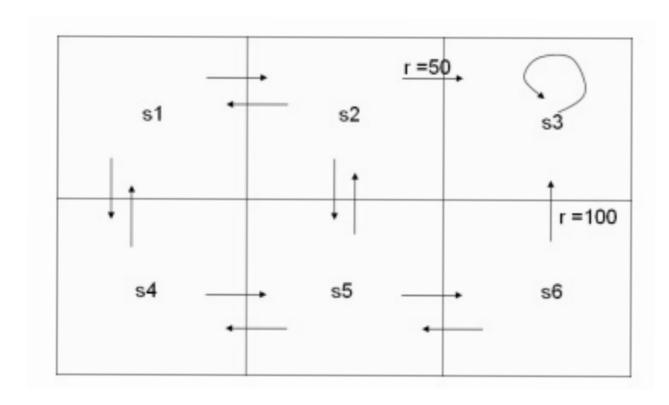
10-701 Final Exam, Fall 2006

9 MDPs and Reinforcement Learning [16pts]

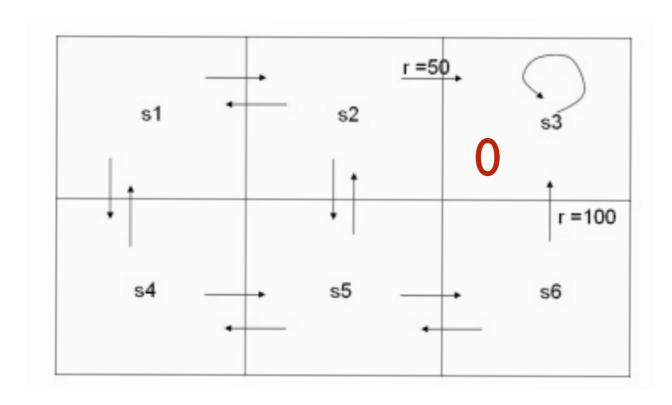
Part A. [10pts]

Consider the following deterministic Markov Decision Process (MDP), describing a simple robot grid world. Notice the values of the *immediate rewards* are written next to transitions. Transitions with no value have an immediate reward of 0. **Assume the discount factor** $\gamma = 0.8$.

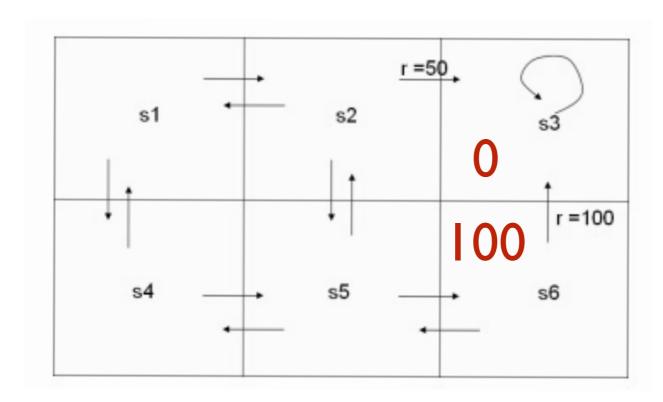




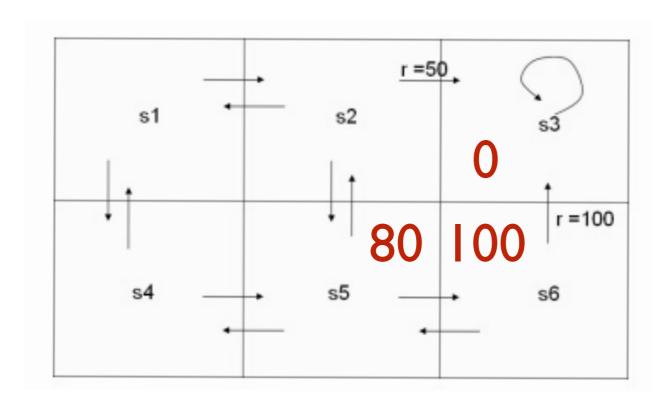
$$\gamma = 0.8$$



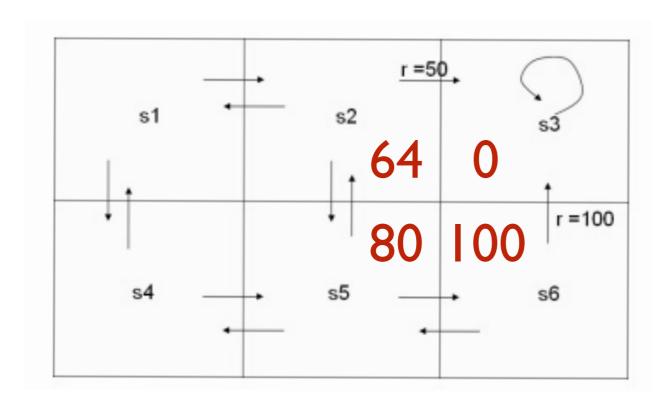
$$\gamma = 0.8$$



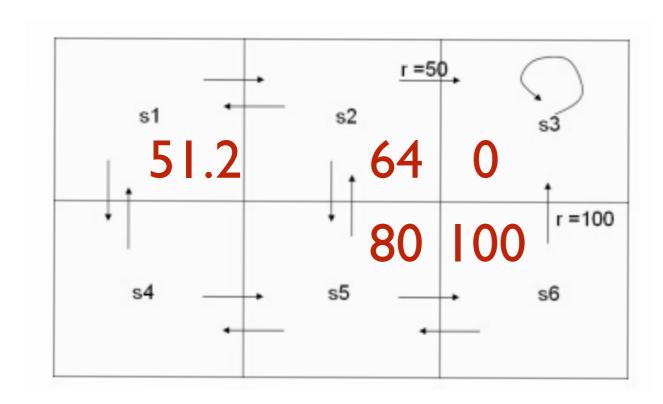
$$y = 0.8$$



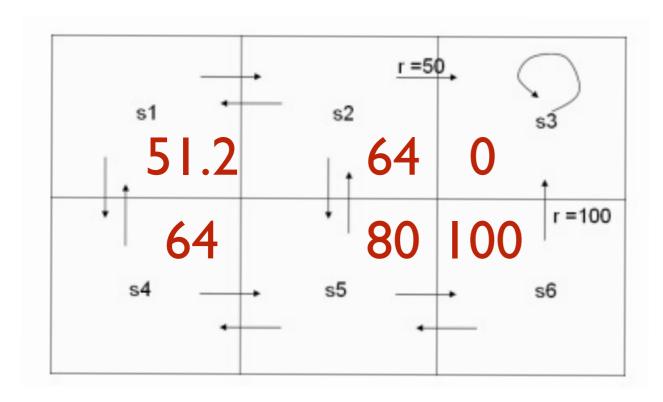
$$y = 0.8$$



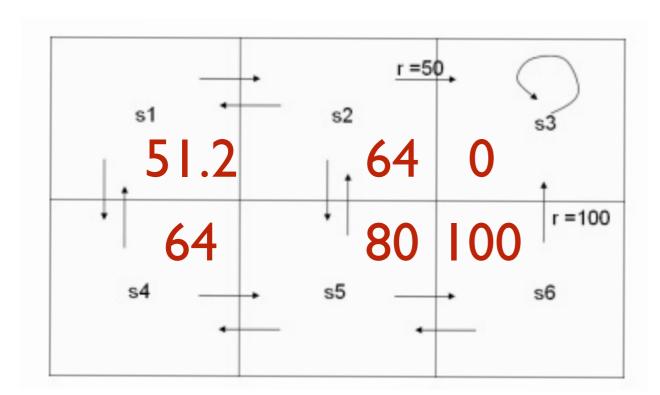
$$y = 0.8$$



$$y = 0.8$$

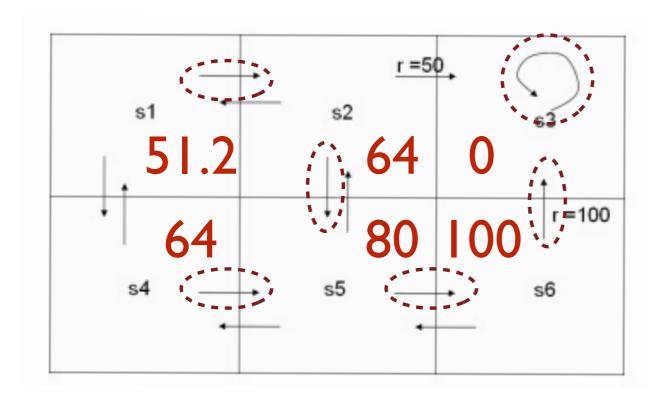


$$\gamma = 0.8$$



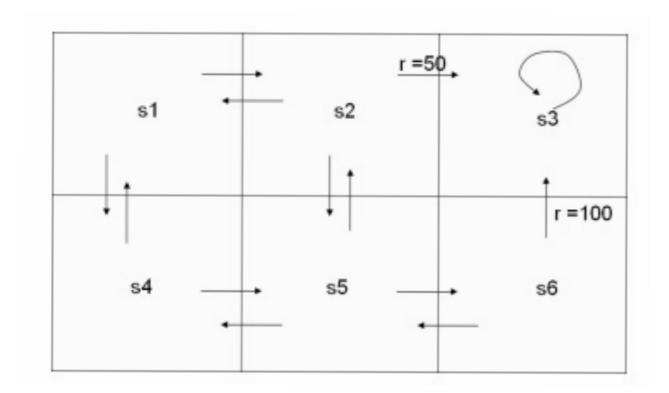
$$\gamma = 0.8$$

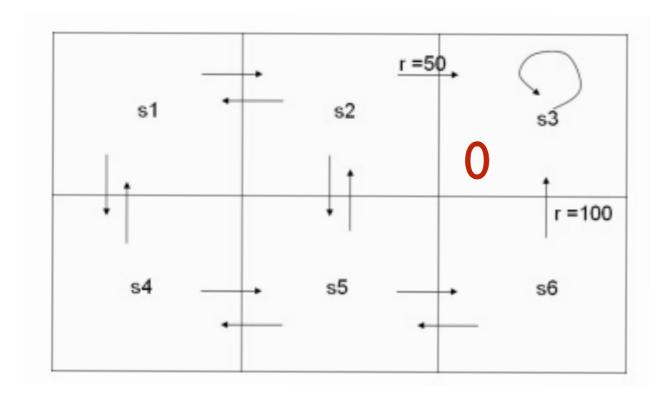
B. Mark the state-action transition arrows that correspond to one optimal pol If there is a tie, always choose the state with the smallest index.

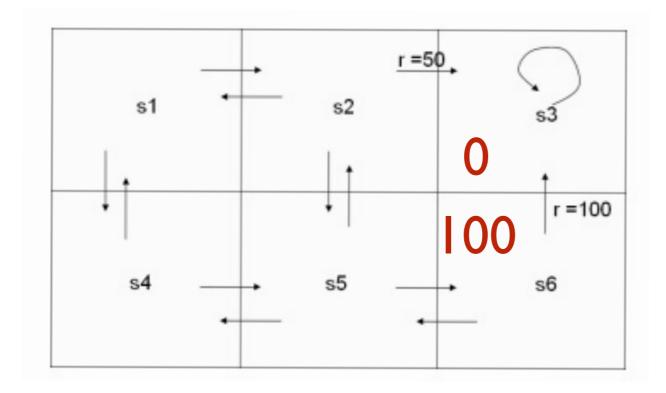


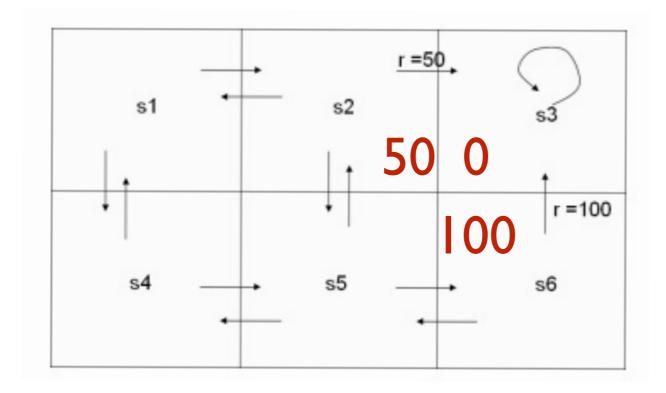
$$\gamma = 0.8$$

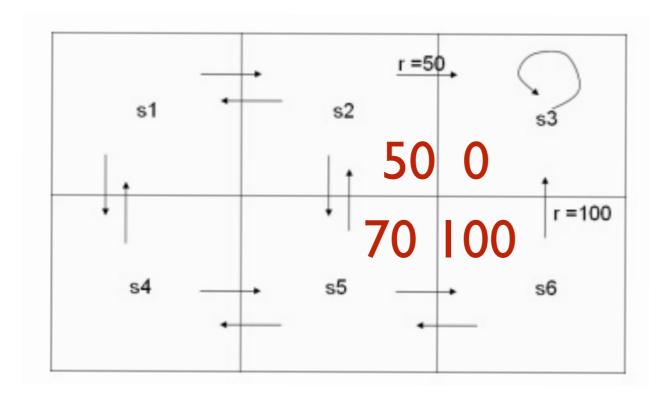
B. Mark the state-action transition arrows that correspond to one optimal pol If there is a tie, always choose the state with the smallest index.

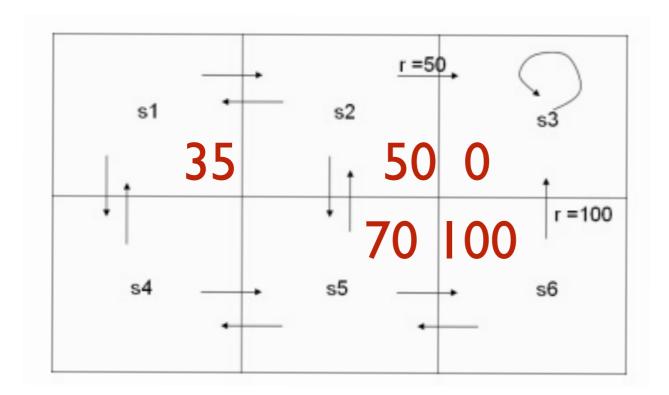


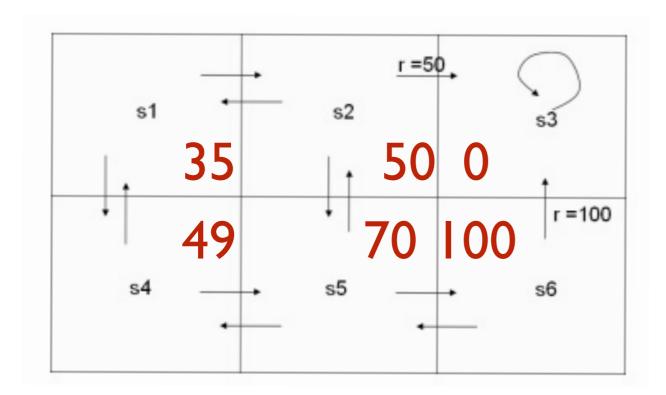


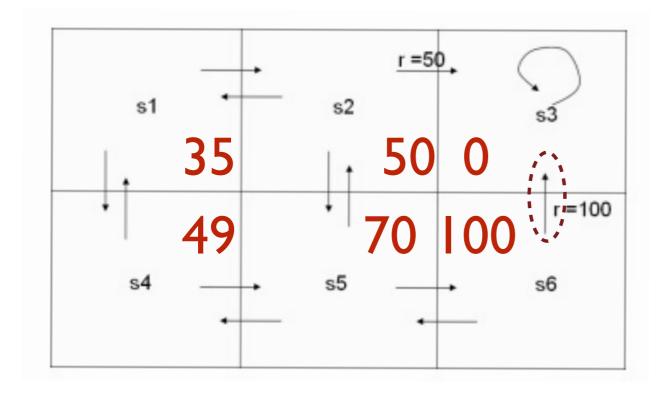


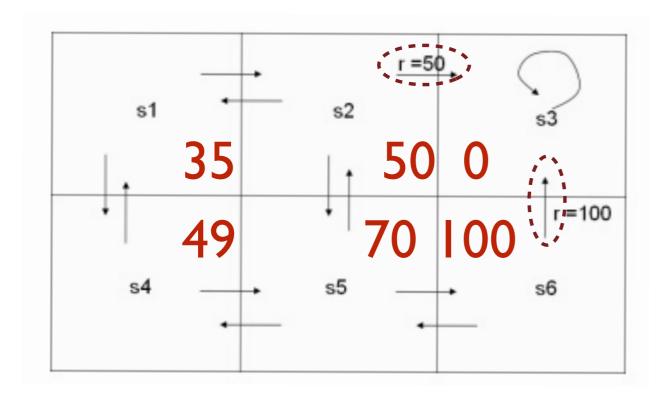


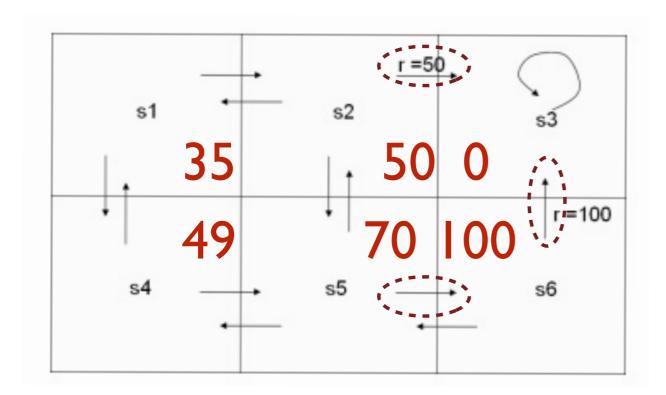


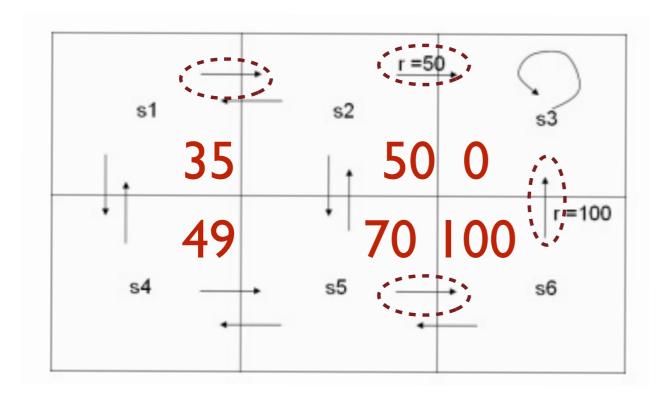


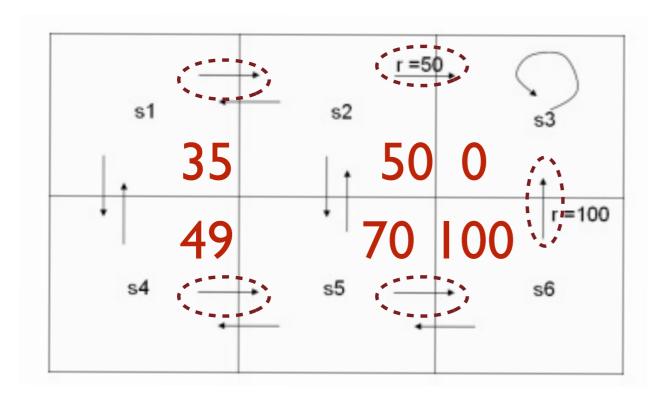


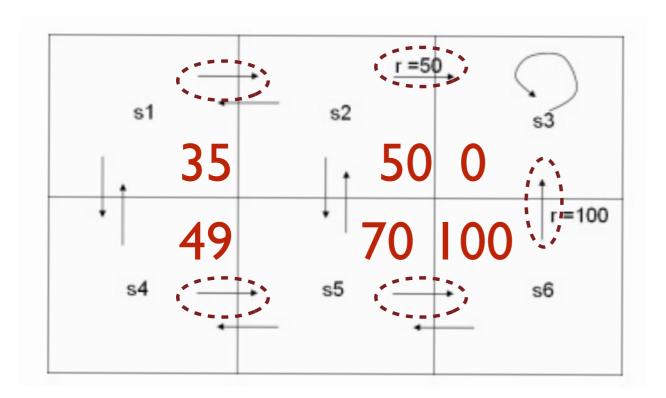








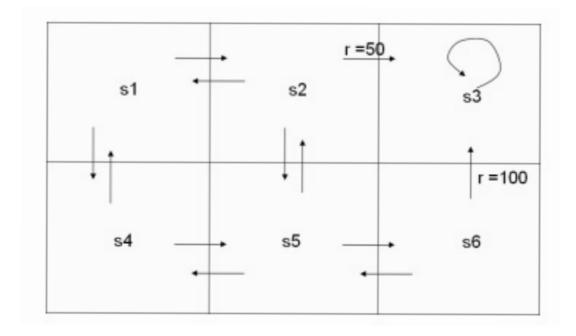




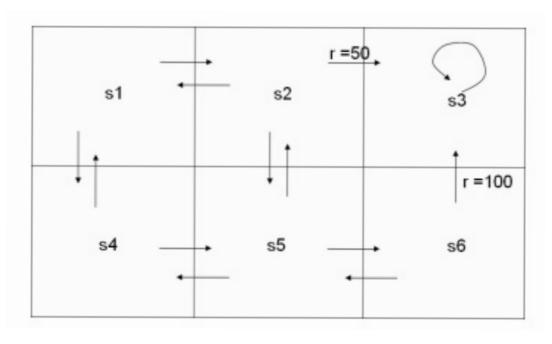
New value for γ : 0.7

Changed policy actions: \$2 -> \$3

A4. How many complete loops (iterations) of value iteration are sufficient to guarantee finding the optimal policy for this MDP? Assume that values are initialized to zero, and that states are considered in an arbitrary order on each iteration.



A4. How many complete loops (iterations) of value iteration are sufficient to guarantee finding the optimal policy for this MDP? Assume that values are initialized to zero, and that states are considered in an arbitrary order on each iteration.



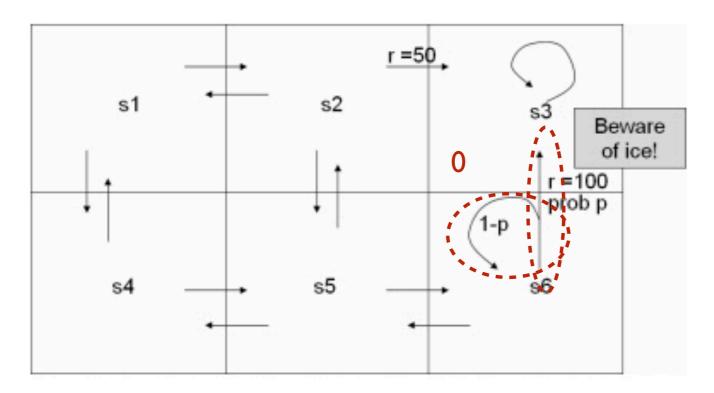
	SI	S2	S3	S4	S5	S6
T=0	0	0	0	0	0	0
T=I	0	50	0	0	0	100
T=2	40	50	0	0	80	100
T=3	40	64	0	64	80	100
T=4	51.2	64	0	64	80	100

A5. Is it possible to change the immediate reward function so that V^* changes but the optimal policy π^* remains unchanged? If yes, give such a change, and describe the resulting change to V. Otherwise, explain in at most 2 sentences why this is impossible.

A5. Is it possible to change the immediate reward function so that V^* changes but the optimal policy π^* remains unchanged? If yes, give such a change, and describe the resulting change to V. Otherwise, explain in at most 2 sentences why this is impossible.

Modify each reward equally, V* will change but the optimal policy will remain the same

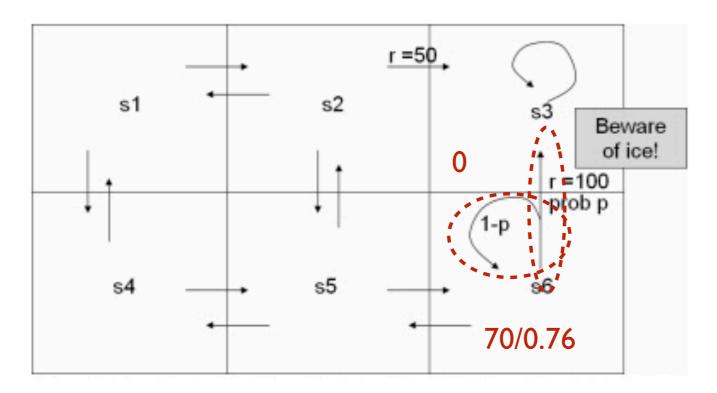
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \\ V_6^* = \frac{70}{0.76} \end{cases}$$

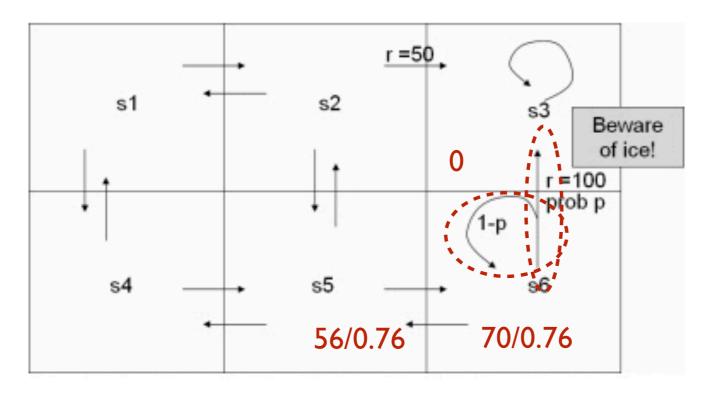
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \end{cases}$$
$$V_6^* = \frac{70}{0.76}$$

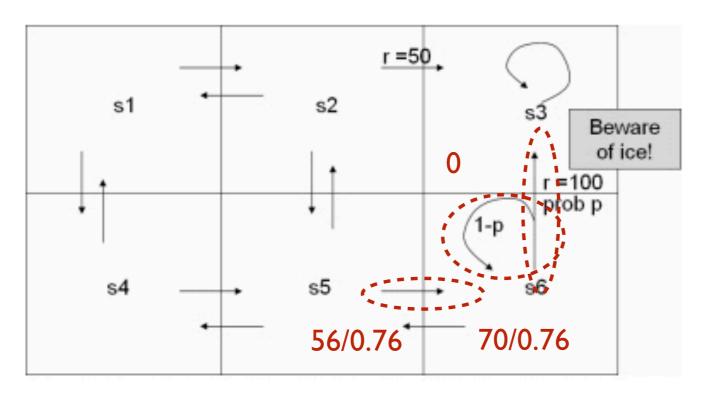
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \end{cases}$$
$$V_6^* = \frac{70}{0.76}$$

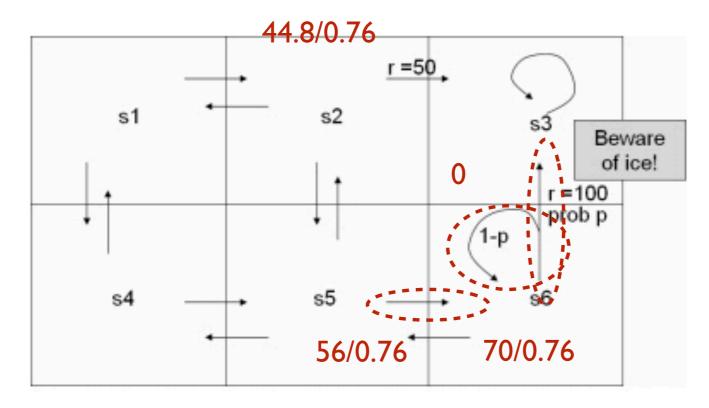
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \\ V_6^* = \frac{70}{0.76} \end{cases}$$

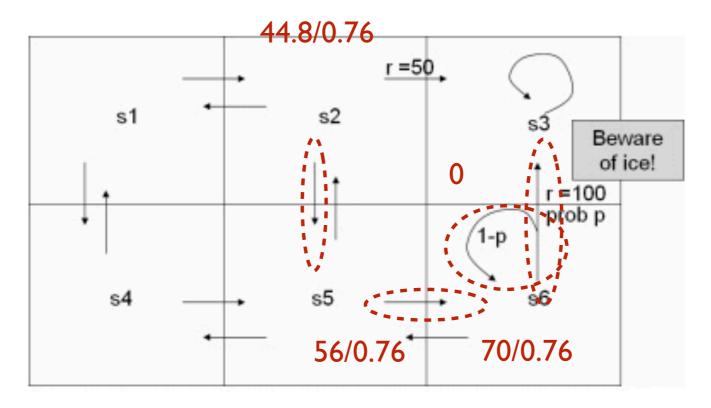
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \end{cases}$$
$$V_6^* = \frac{70}{0.76}$$

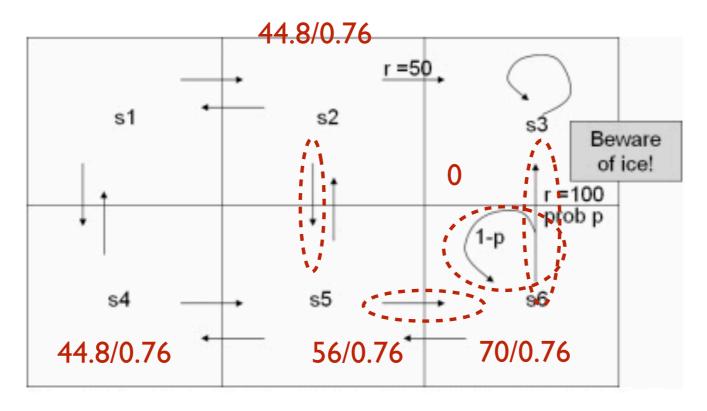
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \end{cases}$$
$$V_6^* = \frac{70}{0.76}$$

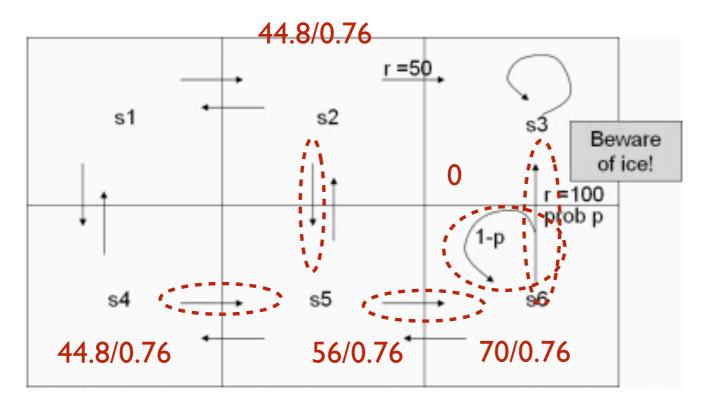
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \end{cases}$$
$$V_6^* = \frac{70}{0.76}$$

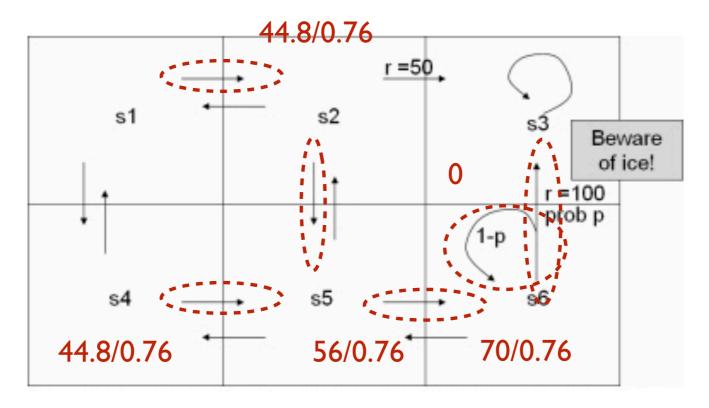
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \\ V_6^* = \frac{70}{0.76} \end{cases}$$

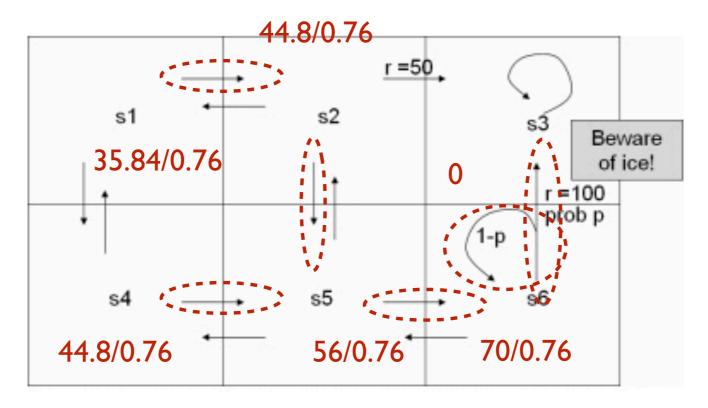
It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.



$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

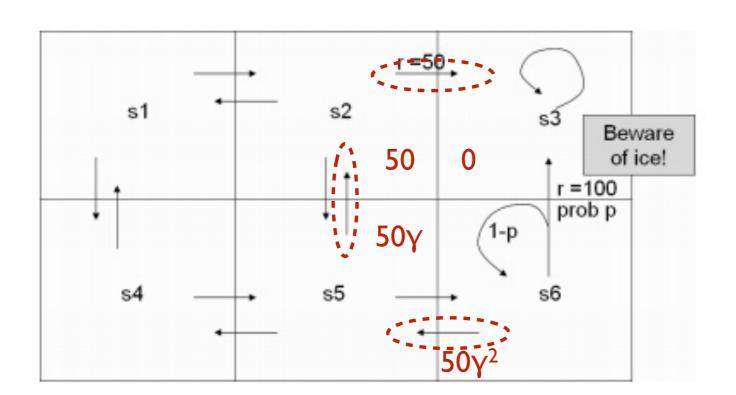
$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \\ V_6^* = \frac{70}{0.76} \end{cases}$$

It is December. Unfortunately for our robot, a patch of ice has appeared in its world, making one of its actions non-deterministic. The resulting MDP is shown below. Note that now the result of the action "go north" from state s6 results in one of two outcomes. With probability p the robot succeeds in transitioning to state s3 and receives immediate reward 100. However, with probability (1-p) it slips on the ice, and remains in state s6 with zero immediate reward. Assume the discount factor $\gamma = 0.8$.

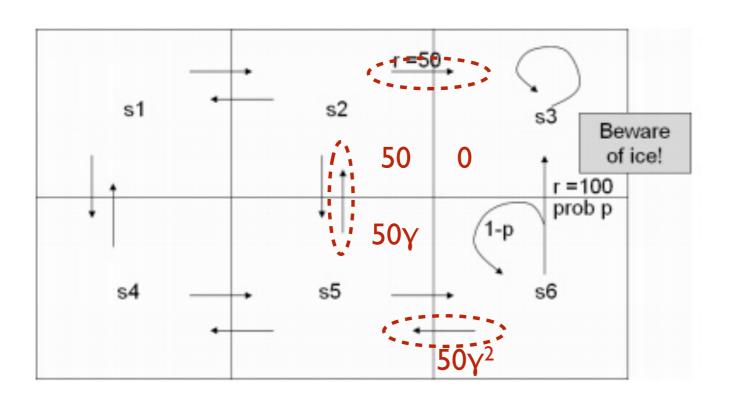


$$V^*(s) = \sum_{s'} P(s'|\pi(s), s) r(s', a, s) + \gamma P(s'|\pi(s), s) V^*(s')$$

$$\begin{cases} V_6^* = 100p + \gamma(1-p)V_6^* \\ V_6^* = 70 + 0.24V_6^* \\ V_6^* = \frac{70}{0.76} \end{cases}$$



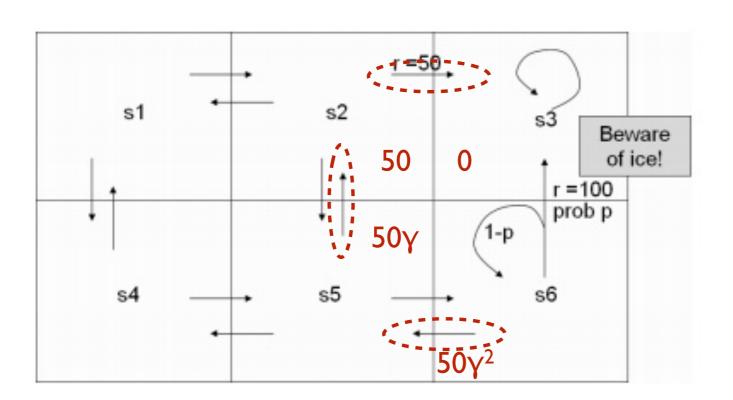
$$V_6^* = 100p + \gamma(1-p)V_6^*$$
$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$



$$V_6^* = 100p + \gamma(1-p)V_6^*$$

$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$

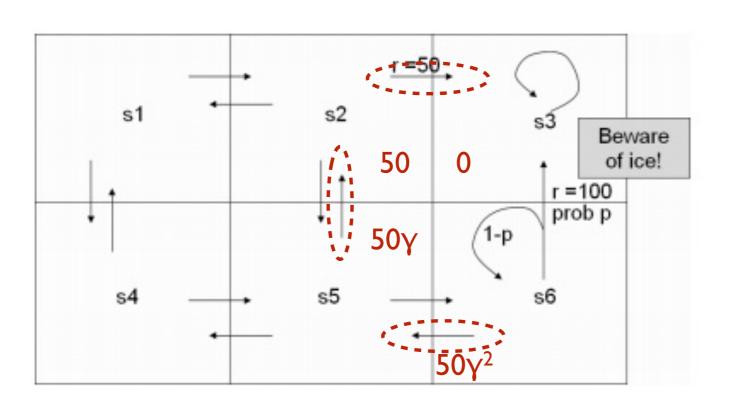
$$50\gamma^2 = V_6^*$$



$$V_6^* = 100p + \gamma(1-p)V_6^*$$
$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$

$$50\gamma^2 = V_6^*$$

$$50\gamma^2 = \frac{100p}{1 - \gamma(1 - p)}$$

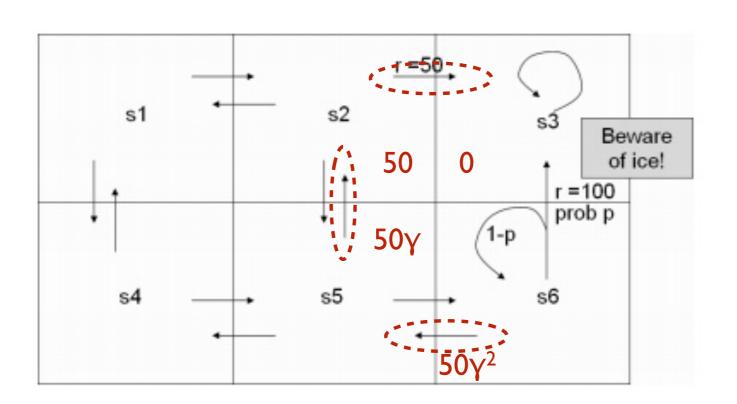


$$V_6^* = 100p + \gamma(1-p)V_6^*$$
$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$

$$50\gamma^2 = V_6^*$$

$$50\gamma^2 = \frac{100p}{1 - \gamma(1 - p)}$$

$$\gamma = 0.8$$



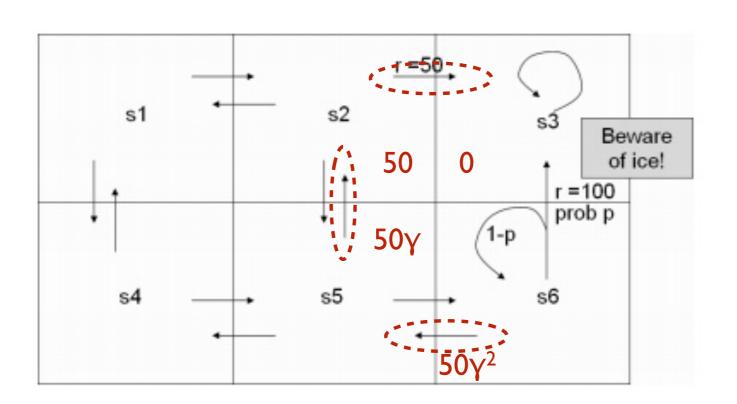
$$V_6^* = 100p + \gamma(1-p)V_6^*$$
$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$

Also...

$$50\gamma^2 = V_6^*$$

$$50\gamma^2 = \frac{100p}{100p}$$

Solve the eq.



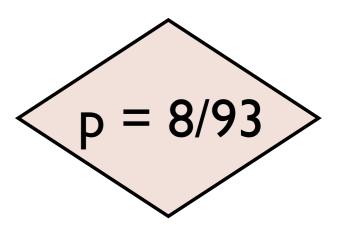
$$V_6^* = 100p + \gamma(1-p)V_6^*$$
$$V_6^* = \frac{100p}{1-\gamma(1-p)}$$

Also...

$$50\gamma^2 = V_6^*$$

$$\frac{100p}{100p}$$

Solve the eq.



QUESTIONS?

REFERENCES

- A good reading is: "Reinforcement Learning: A Survey (1996)" Leslie Pack Kaelbling, Michael L. Littman, Andrew W. Moore. Journal of Artificial Intelligence Research, 4, 237-285. Here is a link where you can download the paper: http://www.autonlab.org/autonweb/14686/version/3/part/5/data/kaelbling-reinforcement.pdf?branch=main&language=en
- Marc Toussaint lecture notes (retrieved from http://userpage.fu-berlin.de/
 mtoussai/notes/markov-decision-processes.pdf