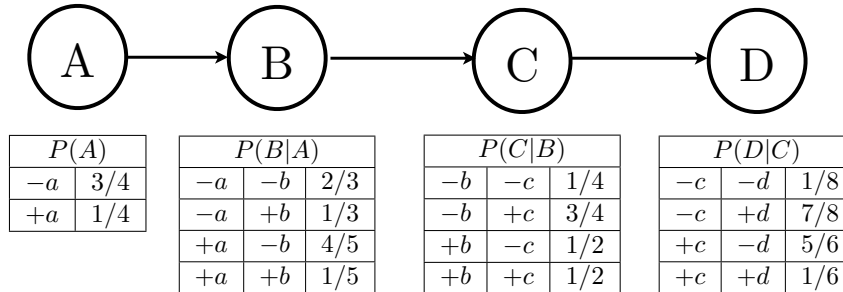


Q1. Bayes' Nets Sampling

Assume the following Bayes' net, and the corresponding distributions over the variables in the Bayes' net:



(a) You are given the following samples:

$+a$	$+b$	$-c$	$-d$	$+a$	$-b$	$-c$	$+d$
$+a$	$-b$	$+c$	$-d$	$+a$	$+b$	$+c$	$-d$
$-a$	$+b$	$+c$	$-d$	$-a$	$+b$	$-c$	$+d$
$-a$	$-b$	$+c$	$-d$	$-a$	$-b$	$+c$	$-d$

(i) Assume that these samples came from performing Prior Sampling, and calculate the sample estimate of  $P(+c)$ .

(ii) Now we will estimate  $P(+c \mid +a, -d)$ . Above, clearly cross out the samples that would **not** be used when doing Rejection Sampling for this task, and write down the sample estimate of  $P(+c \mid +a, -d)$  below.

(b) Using Likelihood Weighting Sampling to estimate  $P(-a \mid +b, -d)$ , the following samples were obtained. Fill in the weight of each sample in the corresponding row.

Sample	Weight
$-a \quad +b \quad +c \quad -d$	_____
$+a \quad +b \quad +c \quad -d$	_____
$+a \quad +b \quad -c \quad -d$	_____
$-a \quad +b \quad -c \quad -d$	_____

(c) From the weighted samples in the previous question, estimate  $P(-a \mid +b, -d)$ .

(d) Which query is better suited for likelihood weighting,  $P(D | A)$  or  $P(A | D)$ ? Justify your answer in one sentence.

(e) Recall that during Gibbs Sampling, samples are generated through an iterative process.

Assume that the only evidence that is available is  $A = +a$ . Clearly fill in the circle(s) of the sequence(s) below that could have been generated by Gibbs Sampling.

Sequence 1				
1 :	$+a$	$-b$	$-c$	$+d$
2 :	$+a$	$-b$	$-c$	$+d$
3 :	$+a$	$-b$	$+c$	$+d$

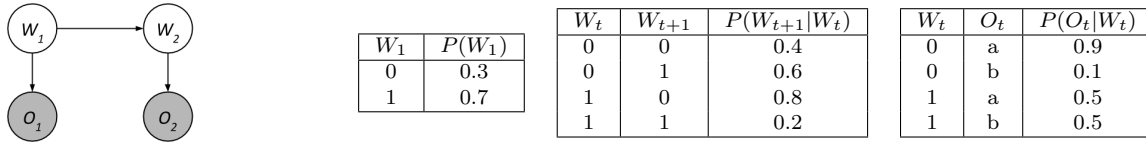
Sequence 2				
1 :	$+a$	$-b$	$-c$	$+d$
2 :	$+a$	$-b$	$-c$	$-d$
3 :	$-a$	$-b$	$-c$	$+d$

Sequence 3				
1 :	$+a$	$-b$	$-c$	$+d$
2 :	$+a$	$-b$	$-c$	$-d$
3 :	$+a$	$+b$	$-c$	$-d$

Sequence 4				
1 :	$+a$	$-b$	$-c$	$+d$
2 :	$+a$	$-b$	$-c$	$-d$
3 :	$+a$	$+b$	$-c$	$+d$

## 2 Particle Filtering

Let's use Particle Filtering to estimate the distribution of  $P(W_2|O_1 = a, O_2 = b)$ . Here's the HMM again.  $O_1$  and  $O_2$  are supposed to be shaded.



We start with two particles representing our distribution for  $W_1$ .

$P_1 : W_1 = 0$

$P_2 : W_1 = 1$

Use the following random numbers to run particle filtering:

[0.22, 0.05, 0.33, 0.20, 0.84, 0.54, 0.79, 0.66, 0.14, 0.96]

(a) **Observe:** Compute the weight of the two particles after evidence  $O_1 = a$ .

(b) **Resample:** Using the random numbers, resample  $P_1$  and  $P_2$  based on the weights.

(c) **Predict:** Sample  $P_1$  and  $P_2$  from applying the time update.

(d) **Update:** Compute the weight of the two particles after evidence  $O_2 = b$ .

(e) **Resample:** Using the random numbers, resample  $P_1$  and  $P_2$  based on the weights.

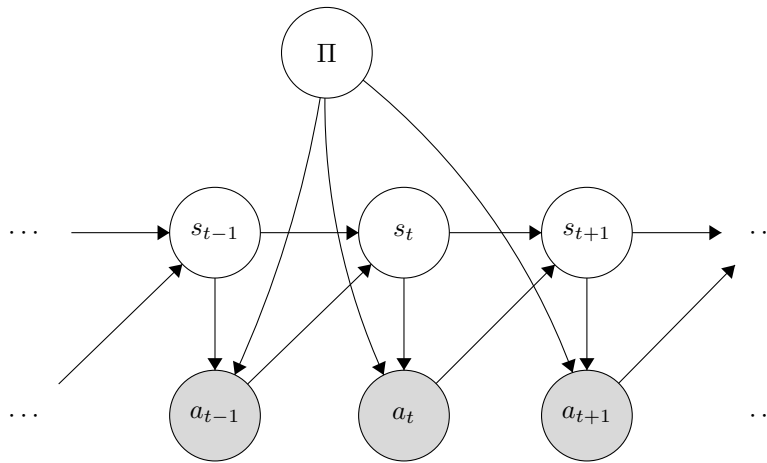
(f) What is our estimated distribution for  $P(W_2|O_1 = a, O_2 = b)$ ?

### 3 [Optional] Particle Filtering Apprenticeship

Consider a modified version of the apprenticeship problem. We are observing an agent's actions in an MDP and are trying to determine which out of a set  $\{\pi_1, \dots, \pi_n\}$  the agent is following. Let the random variable  $\Pi$  take values in that set and represent the policy that the agent is acting under. We consider only *stochastic* policies, so that  $A_t$  is a random variable with a distribution conditioned on  $S_t$  and  $\Pi$ . As in a typical MDP,  $S_t$  is a random variable with a distribution conditioned on  $S_{t-1}$  and  $A_{t-1}$ . The full Bayes net is shown below.

The agent acting in the environment knows what state it is currently in (as is typical in the MDP setting). Unfortunately, however, we, the observer, cannot see the states  $S_t$ . Thus we are forced to use an adapted particle filtering algorithm to solve this problem. Concretely, we will develop an efficient algorithm to estimate  $P(\Pi | a_{1:t})$ .

(a) The Bayes net for part (a) is



(i) Select all of the following that are guaranteed to be true in this model for  $t > 10$ :

- |  |   |
|--|---|
| <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   S_{t-1}$            | <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   \Pi, S_{t-1}$            |
| <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   S_{t-1}, A_{1:t-1}$ | <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   \Pi, S_{t-1}, A_{1:t-1}$ |
| <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   \Pi$                | <input type="checkbox"/> None of the above  |
| <input type="checkbox"/> $S_t \perp\!\!\!\perp S_{t-2}   \Pi, A_{1:t-1}$     |   |

We will compute our estimate for  $P(\Pi | a_{1:t})$  by coming up with a recursive algorithm for computing  $P(\Pi, S_t | a_{1:t})$ . (We can then sum out  $S_t$  to get the desired distribution; in this problem we ignore that step.)

(ii) Write a recursive expression for  $P(\Pi, S_t | a_{1:t})$  in terms of the CPTs in the Bayes net above. Hint: Think of the forward algorithm.

$$P(\Pi, S_t | a_{1:t}) \propto \underline{\hspace{15em}}$$

We now try to adapt particle filtering to approximate this value. Each particle will contain a single state  $s_t$  and a potential policy  $\pi_t$ .

(iii) The following is pseudocode for the body of the loop in our adapted particle filtering algorithm. Fill in the boxes with the correct values so that the algorithm will approximate  $P(\Pi, S_t | a_{1:t})$ .

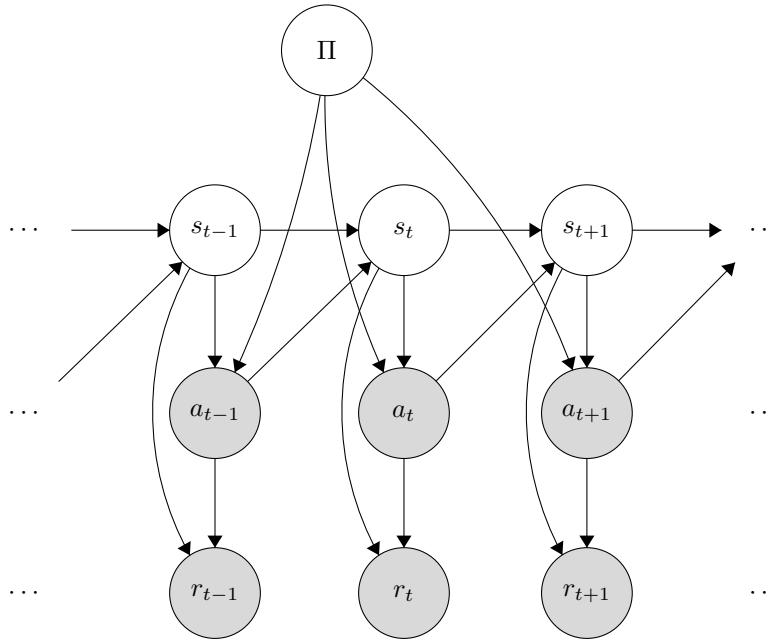
1. Elapse time: for each particle  $(s_t, \pi_t)$ , sample a successor  $s_{t+1}$  from

. The policy  $\pi'$  in the new particle is

2. Incorporate evidence: To each new particle  $(s_{t+1}, \pi')$ , assign weight

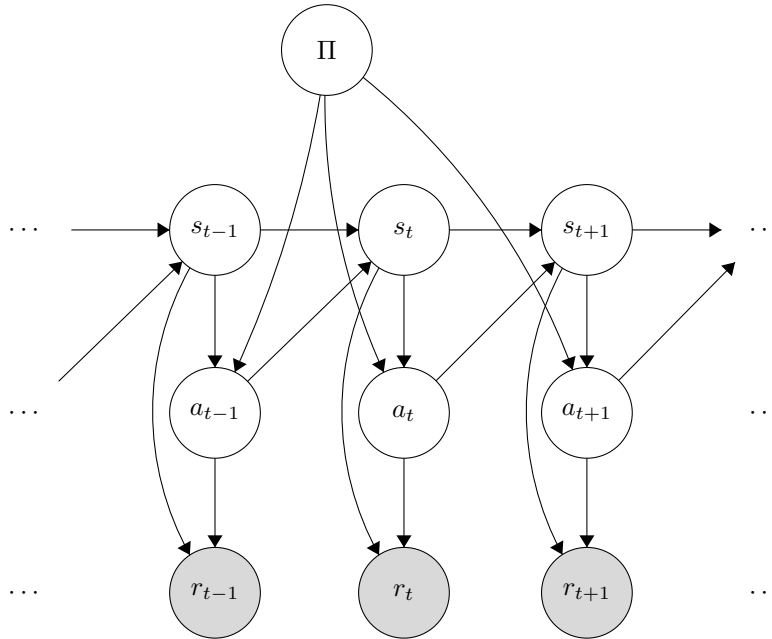
3. Resample particles from the weighted particle distribution.

(b) We now observe the acting agent's actions *and* rewards at each time step (but we still don't know the states). Unlike the MDPs in lecture, here we use a stochastic reward function, so that  $R_t$  is a random variable with a distribution conditioned on  $S_t$  and  $A_t$ . The new Bayes net is given by



Notice that the observed rewards do in fact give useful information since d-separation does not give that  $R_t \perp\!\!\!\perp \Pi | A_{1:t}$ . Give an active path connecting  $R_t$  and  $\Pi$  when  $A_{1:t}$  are observed. Your answer should be an ordered list of nodes in the graph, for example “ $S_t, S_{t+1}, A_t, \Pi, A_{t-1}, R_{t-1}$ ”.

(c) We now observe *only* the sequence of rewards and no longer observe the sequence of actions. The new Bayes net is:



We will compute our estimate for  $P(\Pi \mid r_{1:t})$  by coming up with a recursive algorithm for computing  $P(\Pi, S_t, A_t \mid r_{1:t})$ . (We can then sum out  $S_t$  and  $A_t$  to get the desired distribution; in this problem we ignore that step.)

(i) Write a recursive expression for  $P(\Pi, S_t, A_t \mid r_{1:t})$  in terms of the CPTs in the Bayes net above.

$$P(\Pi, S_t, A_t \mid r_{1:t}) \propto \underline{\hspace{15em}}$$

We now try to adapt particle filtering to approximate this value. Each particle will contain a single state  $s_t$ , a single action  $a_t$ , and a potential policy  $\pi_i$ .

(ii) The following is pseudocode for the body of the loop in our adapted particle filtering algorithm. Fill in the boxes with the correct values so that the algorithm will approximate  $P(\Pi, S_t, A_t \mid r_{1:t})$ .

1. Elapse time: for each particle  $(s_t, a_t, \pi_i)$ , sample a successor state  $s_{t+1}$  from

. Then, sample a successor action  $a_{t+1}$  from

. The policy  $\pi'$  in the new particle is

2. Incorporate evidence: To each new particle  $(s_{t+1}, a_{t+1}, \pi')$ , assign weight

3. Resample particles from the weighted particle distribution.