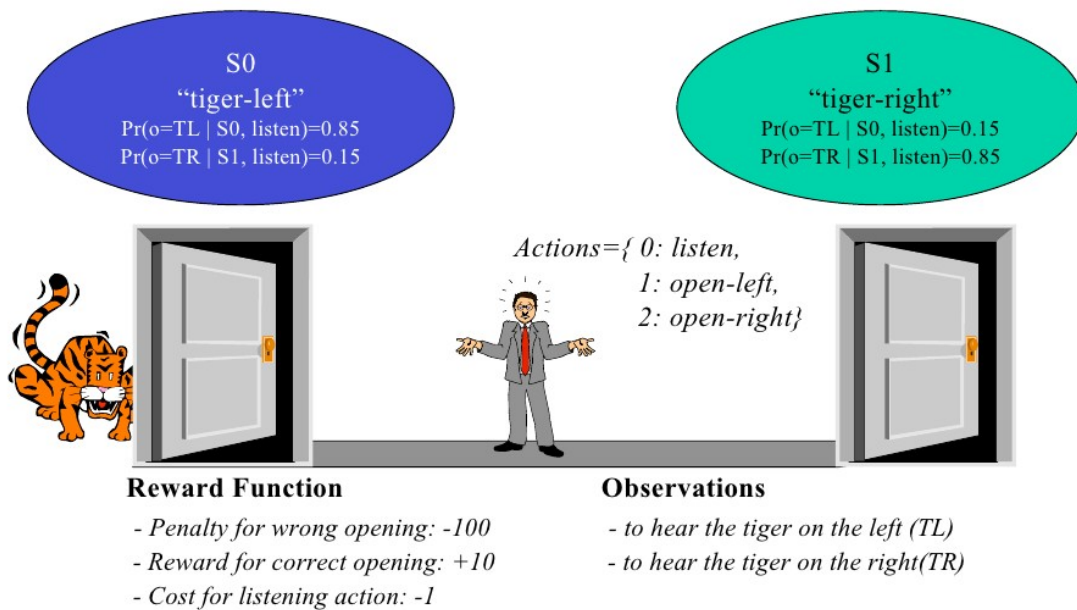# Reinforcement Learning: Tutorial 5 (week from 3. 3. 2014)

1. How can particle filters be used in the context of robot localization?

2. The "art" of importance sampling: We are sampling $P(x)$, which may be not cover the interesting aspect of the game. It is already interesting to consider sampling $X/L(x)$ based on the distribution $P(x)L(x)$. Why? How could it be useful in RL? Similarly, we can consider this scheme (see on the right ->) Why should one want to do this? What happens for $P(x)=W(x)$?

1. Algorithm **Discrete_Bayes_filter**( $Bel(x),d$ )
2. $\eta=0$
3. If $d$ is a perceptual data item $z$ then
4.    For all $x$ do
5.       $Bel'(x) = P(z\,|\,x)Bel(x)$
6.       $\eta = \eta + Bel'(x)$
7.    For all $x$ do
8.       $Bel'(x) = \eta^{-1}Bel'(x)$
9. Else if $d$ is an action data item $u$ then
10.    For all $x$ do
11.       $Bel'(x) = \sum_{x'} P(x\,|\,u,x')\,Bel(x')$
12. Return $Bel'(x)$

$$\overline{A} = \frac{\sum_x P(x)\,A(x)\,/\,W(x)}{\sum_x P(x)\,/\,W(x)}\ .$$

3. How are POMDPs and Hidden Markov Models (HMMs) related? Would a Viterbi algorithm be useful in POMDPs?

4. Discuss the tiger problem (from: Dr. Stephan Timmer "Introduction to POMDPs")



**S0**
"tiger-left"
Pr(o=TL | S0, listen)=0.85
Pr(o=TR | S1, listen)=0.15

**S1**
"tiger-right"
Pr(o=TL | S0, listen)=0.15
Pr(o=TR | S1, listen)=0.85

Actions={ 0: listen,
1: open-left,
2: open-right}

**Reward Function**
- Penalty for wrong opening: -100
- Reward for correct opening: +10
- Cost for listening action: -1

**Observations**
- to hear the tiger on the left (TL)
- to hear the tiger on the right(TR)

```
# This is the tiger problem of AAAI paper fame in the new POMDP
# format. This format is still experimental and subject to change

discount: 0.75
values: reward
```
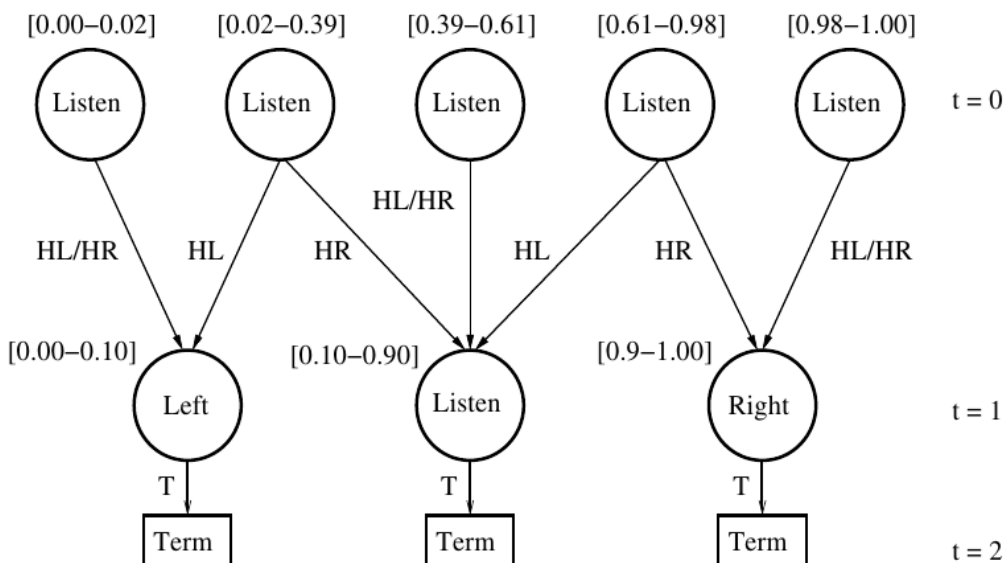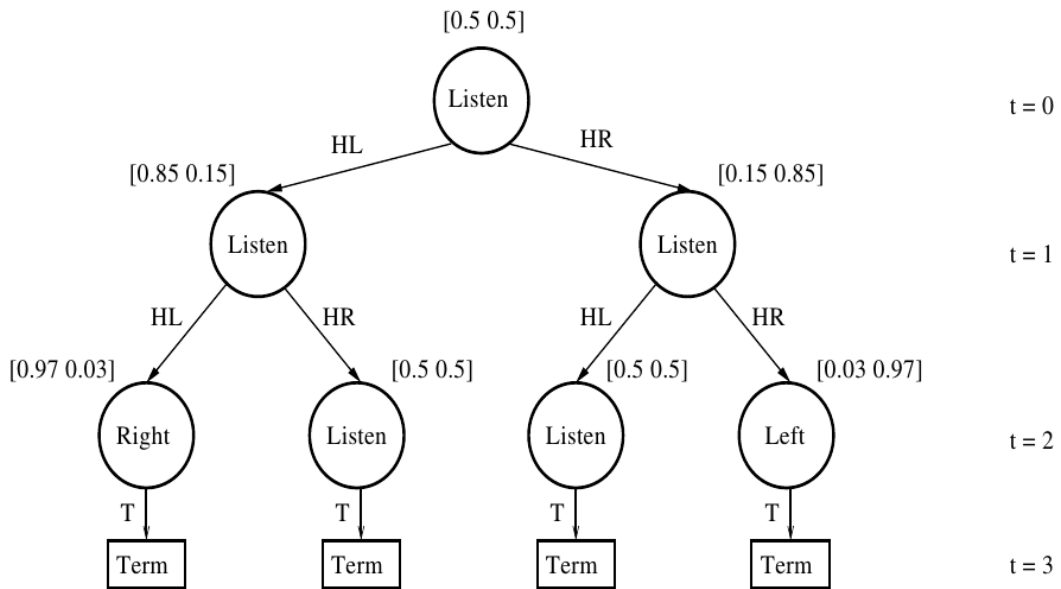
```
states: tiger-left tiger-right
actions: listen, open-left, open-right
observations: tiger-left, tiger-right

Transitions:
listen -> identity
open-left -> uniform
open-right -> uniform

Observations
listen (in either state):      0.85 0.15
                               0.15 0.85
open-left: uniform
open-right: uniform

Rewards:
R:listen : * : * : * -1
R:open-left : tiger-left : * : * -100
R:open-left : tiger-right : * : * 10
R:open-right : tiger-left : * : * 10
R:open-right : tiger-right : * : * -100
```

5. Consider the application to a POMDP to the problem of controlling several elevators problem. For what definition of states does any uncertainty arise? Discuss the advantage of a POMDP over a state abstraction (that does not distinguish between states that can be confused). Compare to the original Barto&Crites approach (see final slides of the lecture RL09). Can the design of the elevator operation be changed such that this uncertainty is removed/reduced?

6. Recall the discussion of "afterstates" from a previous tutorial. Afterstates are an option to include the reaction of an opponent into the own policy. Under what conditions would it make sense to reformulate the problem as POMDPs.

7. Consider a robot moving down a hallway as a 1D problem with states being sections of the track of a length of 1m. The robot's speed is 1m/s +/- 0.1m/s (assume a uniform distribution of deviations). Discuss the belief propagation in standard POMDP vs. the corresponding effects in an augmented MDP or in QMDP. Think of a navigation task which is then to be solved by either of these methods.

8. Assume a robot moving in a dark environment where information is available only from touch sensors. The robot learns to move successfully using a POMDP. Now the lights are switched on and the robot can use again its excellent visual system. How can it use the information from POMDP for initialising a simpler RL method

9. Have a look at a review paper such as Anthony R. Cassandra (1998) A Survey of POMDP Applications. Discuss set-up, advantages and limitations of POMDP in the mentioned application problems (or do this simply for any of the examples above).