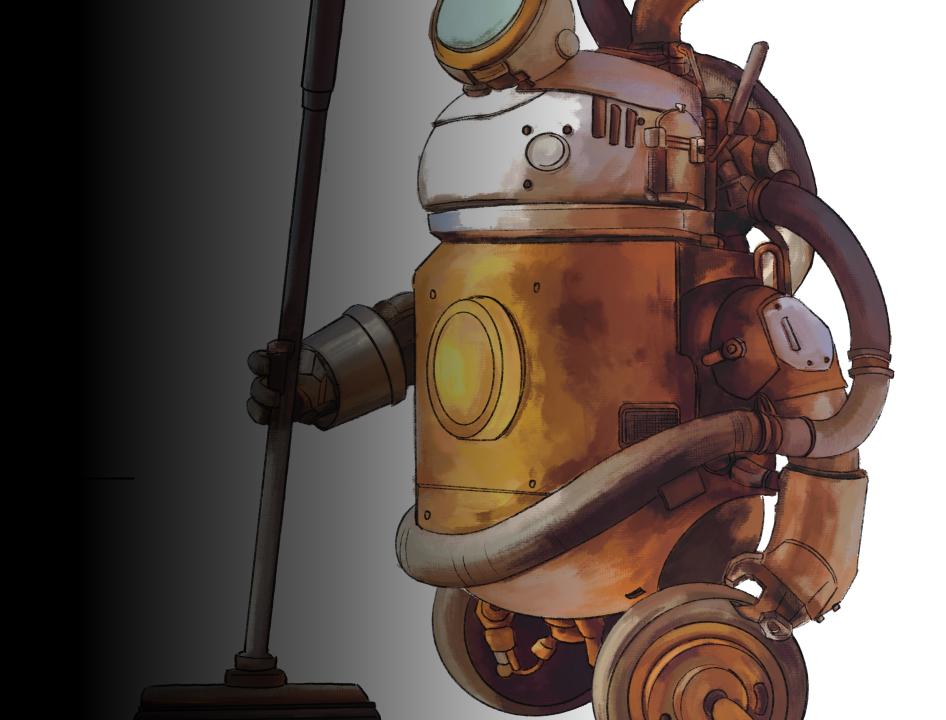
# CS 3630, Fall 2025

Lecture 10:

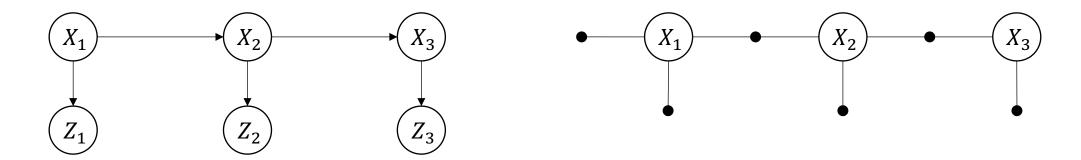
**Markov Decision Processes** 



Lecture 9 Recap



## Factor Graphs



Measurements are given – get rid of them!

$$P(X|Z) \propto P(X_1)L(X_1; z_1)P(X_2|X_1)L(X_2; z_2)P(X_3|X_2)L(X_3; z_3)$$

• This becomes:

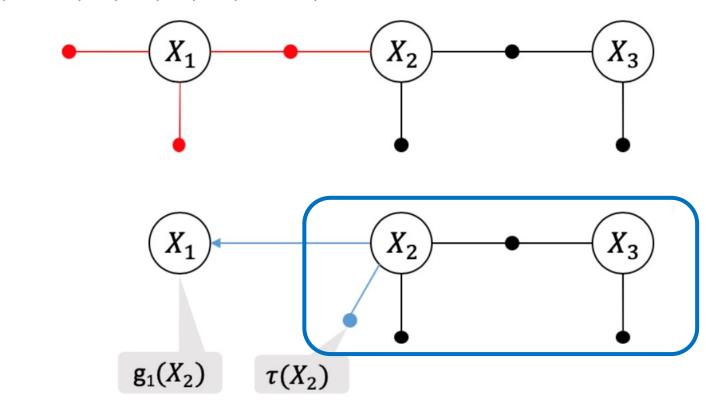
$$\phi(\mathcal{X}) = \phi_1(X_1)\phi_2(X_1)\phi_3(X_1, X_2)\phi_4(X_2)\phi_5(X_2, X_3)\phi_6(X_3)$$

Each factor defines a function  $\phi$  which is a function only of its (non-factor node) neighbors.

## MPE via max-product

• Eliminate one variable at a time by forming product, then max:

$$\phi(X_1,X_2) = \phi_1(X_1)\phi_2(X_1)\phi_3(X_1,X_2)$$
 .

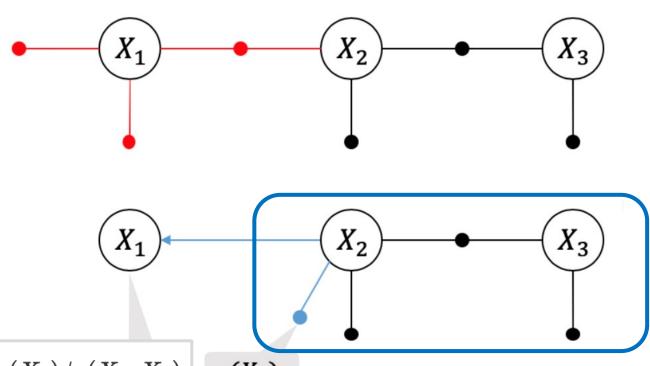


$$g_1(X_2) = rg \max_{x_1} \phi(x_1, X_2) \qquad au(X_2) = \max_{x_1} \phi(x_1, X_2)$$

## Posterior via sum-product:

• Eliminate one variable at a time by forming product, then sum:

$$\phi(X_1,X_2) = \phi_1(X_1)\phi_2(X_1)\phi_3(X_1,X_2)$$



$$P(X_1|X_2) = rac{\phi_1(X_1)\phi_2(X_1)\phi_3(X_1,X_2)}{ au(X_2)}.$$

$$\tau(X_2)$$

$$au(X_2) \doteq \sum_{X_1} \phi_1(X_1) \phi_2(X_1) \phi_3(X_1, X_2)$$

## Markov Decision Processes

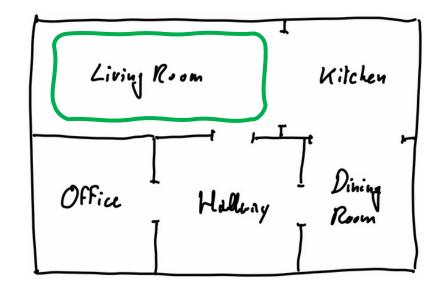
- Planning is the process of choosing which actions to perform.
- In order to plan effectively, we need quantitative criteria to evaluate actions and their effects.
- MDPs include a reward function that characterizes the immediate benefit of applying an action.
- Policies describe how to act in a given state.
- The value function characterizes the long-term benefits of a policy.
- We assume that the robot is able to *know* its current state with certainty.
- > We will see how to define reward functions and use these to compute optimal policies for MDPs.

#### Reward Functions

 Most general form depends on current state, action, and next state:

$$R: \mathcal{X} \times \mathcal{A} \times \mathcal{X} \to \mathbb{R}$$

 In our example, we just care about where we end up after taking an action:



```
def reward_function(state:str, action:str, next_state:str):
    """Reward that returns 10 upon entering the living room."""
    return 10.0 if next_state == "Living Room" else 0.0

print(reward_function("Kitchen", "L", "Living Room"))
print(reward_function("Kitchen", "L", "Kitchen"))
```

```
10.0
0.0
```

## **Expected Reward**

 A greedy way to act would be to calculate the immediate expected reward for every possible action:

$$\overline{R}(x,a) = E[R(x,a,X')]$$

• Since we know the transition probabilities, we can easily compute this:

$$ar{R}(x,a) \doteq E[R(x,a,X')] = \sum_{x'} P(x'|x,a)R(x,a,x')$$

• We then have a simple greedy planning algorithm:

$$a^* = rg \max_{a \in \mathcal{A}} E[R(X_t, a, X_{t+t})]$$

## Example

Expected reward (Kitchen, L) = 8.0Expected reward (Kitchen, R) = 0.0Expected reward (Kitchen, U) = 0.0Expected reward (Kitchen, D) = 0.0

 The expected immediate reward for all four actions in the Kitchen:

```
x = vacuum.rooms.index("Kitchen")
for a in range(4):
    print(f"Expected reward ({vacuum.rooms[x]}, {vacuum.action_space[a]}) = {T[x,a] @ R[x,a]}")

$\square 0.9s$
```

Kitchen

- Hence, when in the kitchen, always do L!
- This is a greedy planning algorithm

## Utility

$$U: \mathcal{A}^n \times \mathcal{X}^{n+1} \to \mathbb{R}$$

$$U(a_1, \dots, a_n, x_1, \dots x_{n+1}) = R(x_1, a_1, x_2) + \gamma R(x_2, a_2, x_3) + \dots + \gamma^{n-1} R(x_n, a_n, x_{n+1})$$

- Because actions are uncertain, let's look further into the future!
- Introduce a discount factor  $\gamma$  to
  - still bias towards more immediate payoff;
  - allow infinite time horizons:

$$U(a_{1, \dots, x_{1}, \dots}) = \sum_{i=1}^{\infty} \gamma^{i-1} R(x_{i}, a_{i}, x_{i+1})$$



## **Expected Utility**

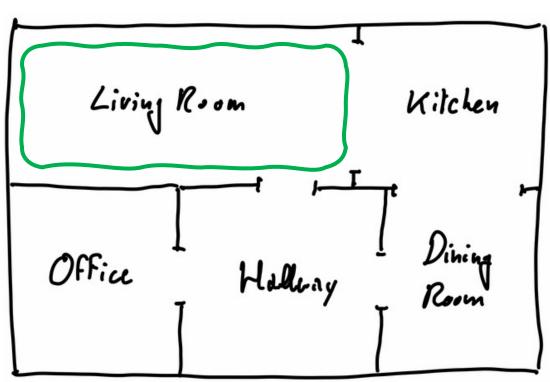
$$E[U(a_1, ..., a_n, x_1, X_2, ... X_{n+1})] = E[R(x_1, a_1, X_2) + \gamma R(X_2, a_2, X_3) + \cdots + \gamma^{n-1} R(X_n, a_n, X_{n+1})]$$

- Again, we can compute the expectation to choose between finite horizon plans
- For n=3, we have  $4^3 = 64$  possible plans, and for each plan we must evaluate  $5^4 = 625$  possible state sequences
- An approximate algorithm to evaluate a given plan:
  - Simulate multiple rollouts
  - Average the result
- Still expensive, only practical for short horizon plans...

### Policies $\pi: \mathcal{X} \to \mathcal{A}$

def reward\_function(state:str, action:str, next\_state:str):
 """Reward that returns 10 upon entering the living room."""
 return 10.0 if next\_state == "Living Room" else 0.0

- Because actions are non-deterministic, fixed plans are brittle and prone to failure.
- Better to have a state-dependent plan
- A policy  $\pi(X)$  is a function that specifies which action to take in each state.
- Let us come up with a policy together:
  - $\pi(L) =$
  - $\pi(K) =$
  - $\pi(0) =$
  - $\pi(H) =$
  - $\pi(D) =$



## The Value Function for a Policy

Recall the Expected Utility

$$\overline{U}(a_1 ... a_n, x_1) = E\left[\sum_{i=1}^n \gamma^{i-1} R(X_i, a_i, X_{i+1})\right]$$

For a policy, we can define this similarly:

$$\overline{U}(\pi, n, x_1) \doteq E\left[R(x_1, \pi(x_1), X_2) + \gamma R(X_2, \pi(X_2), X_3) + \dots + \gamma^2 R(X_n, \pi(X_n), X_n)\right]$$

• Can be extended to infinite horizon policy, defining the value function:

$$V^{\pi}(x_1) \doteq E\left[R(x_1, \pi(x_1), X_2) + \gamma R(X_2, \pi(X_2), X_3) + \gamma^2 R(X_3, \pi(X_3), X_4) + \cdots\right]$$

• Of course, above holds for arbitrary  $x_t$ , not just  $x_1$ .

#### Recursive Definition of $V^{\pi}$

$$V^{\pi}(x_1) = E[R(x_1, \pi(x_1), X_2) + \gamma R(X_2, \pi(X_2), X_3) + \gamma^2 R(X_3, \pi(X_3), X_4) + \dots]$$

$$V^{\pi}(x_1) = \sum_{x_2} P(x_2 | x_1, \pi(x_1)) \{ R(x_1, \pi(x_1), x_2) + \gamma E[R(x_2, \pi(x_2), X_3) + \gamma R(X_3, \pi(X_3), X_4) + \dots] \}$$

$$V^{\pi}(x_1) = \sum_{x_2} P(x_2 | x_1, \pi(x_1)) \{ R(x_1, \pi(x_1), x_2) + \gamma V^{\pi}(x_2) \}$$

$$V^{\pi}(x_1) = \sum_{x_2} P(x_2 | x_1, \pi(x_1)) R(x_1, \pi(x_1), x_2) + \gamma \sum_{x_2} P(x_2 | x_1, \pi(x_1)) V^{\pi}(x_2)$$

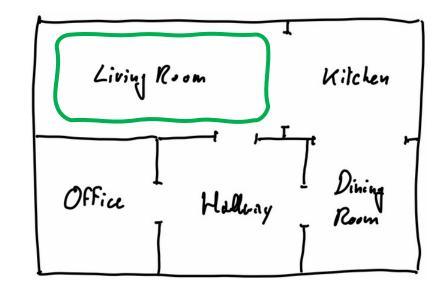
$$V^{\pi}(x) = \bar{R}(x, \pi(x)) + \gamma \sum_{x'} P(x'|x, \pi(x)) V^{\pi}(x')$$

## Exact Computation for $V^{\pi}$

• Because we have a finite set of states, we get 5 linear equations in 5 unknowns  $V^{\pi}(x)$ :

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

- Can be solved efficiently with np.linalg.solve
- Example in book:



```
reasonable_policy = [UP, LEFT, RIGHT, UP, LEFT]
                                                                  V(reasonable_policy):
                                      [[10.]
                                                                    Living Room: 100.00
                                       [ 8.]
[-0.72 \quad 0.82 \quad -0.
                                                                    Kitchen
                                                                                 : 97.56
      -0. 0.82 -0.72 -0. ]
                                       [ 0.]
                                                                    Office
                                                                                 : 85.66
                                       [ 8.]
[-0.72 - 0. -0. 0.82 - 0.]
                                                                    Hallway
                                                                                 : 97.56
                                       [ 0.]]
            -0. -0.72 0.82]]
                                                                    Dining Room: 85.66
```

## Policy Iteration

Start with a random policy  $\pi^0$ , and repeat until convergence:

- 1. Compute the value function  $V^{\pi^k}$
- 2. Improve the policy for each state x using the update rule:

$$\pi^{k+1}(x) \leftarrow \arg\max_{a} \left\{ \overline{R}(x,a) + \gamma \sum_{x'} P(x'|x,a) \right\} V^{\pi^{k}}(x')$$



```
always_right = [RIGHT, RIGHT, RIGHT, RIGHT]

$\square$ 0.7s
```

```
optimal_policy, optimal_value_function = policy_iteration(always_right)
print([vacuum.action_space[a] for a in optimal_policy])

    0.7s

['L', 'L', 'R', 'U', 'U']
```

## Optimal Value Function

The optimal value function is the one corresponding to the optimal policy:

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

$$= \max_{\pi} \left\{ \bar{R}(x, \pi(x)) + \gamma \sum_{x'} P(x'|x, \pi(x)) V^{\pi}(x') \right\}$$

$$= \max_{a} \left\{ \bar{R}(x, a) + \gamma \sum_{x'} P(x'|x, a)) V^{*}(x') \right\}$$

#### The Bellman equation:

$$\overline{V^*(x)} = \max_{a} \left\{ \overline{R}(x,a) + \gamma \sum_{x'} P(x'|x,a)) V^*(x') \right\}$$

#### Value Iteration

Start with a random value function  $V^0$ , and repeat until convergence:

• Improve the value function  $V^k$  using the update rule:

$$V^{k+1}(x) \leftarrow \max_{a} \left\{ \overline{R}(x,a) + \gamma \sum_{x'} P(x'|x,a) \right\}$$

```
V_k = np.full((5,), 100)
for k in range(10):
    Q_k = np.sum(T * (R + 0.9 * V_k), axis=2) # 5 x 4
    V_k = np.max(Q_k, axis=1) # max over actions
    print(np.round(V_k,2))
```

```
[100. 98. 90. 98. 90.]

[100. 97.64 86.76 97.64 86.76]

[100. 97.58 85.92 97.58 85.92]

[100. 97.56 85.72 97.56 85.72]
```

## **Optimal Policy**

Given the  $V^*(x)$ , computing the optimal policy is a straightforward optimization:

$$\pi^*(x) = \arg\max_{a} \left\{ \overline{R}(x,a) + \gamma \sum_{x'} P(x'|x,a)) V^*(x') \right\}$$

For convenience, we define the  $Q^*$  function as

$$Q^*(x,a) = \overline{R}(x,a) + \gamma \sum_{x'} P(x'|x,a))V^{\pi}(x')$$

and we can write the optimal policy as:

$$\pi^*(x) = \arg\max_{a} Q^*(x, a)$$

The Q function plays a role in reinforcement learning, to be continued...