

|               |               |    |
|---------------|---------------|----|
| 06            | 02            | 06 |
| 13            | 03            | 07 |
| 08            | 04            | 22 |
| 09            | 05            | 23 |
| <del>10</del> | <del>14</del> |    |
| <del>11</del> |               |    |
| <del>12</del> | 15            |    |
|               | 13            |    |

18  
19  
20  
21

# oMDP

Finite-stage v.s. ongoing problems  
 infinite horizon: the process may go on forever  
 indefinite horizon: the agent will eventually stop, but it does not know when it will stop.  
 Utility at the end v.s. a sequence of rewards

## -An MDP Needs:

S: set of states  
 A: set of actions

transition probabilities:  $P(s'|s, a)$

Reward function:  $R(s, a, s')$

## -Choice of Reward Function.

$R(s)$ : the reward of entering state  $s$ .

- Total reward  $R(S_0) + R(S_1) + R(S_2) + \dots$ 
  - If the sum is infinite, we cannot compare two policies of an MDP.
- Average reward  $\lim_{k \rightarrow \infty} \frac{1}{k} (R(S_0) + R(S_1) + R(S_2) + \dots)$ 
  - If the total reward is finite, the average reward is zero.
- Discounted reward  $R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots$ 
  - where  $0 < \gamma < 1$  is the discount factor.
  - We prefer getting the same reward sooner rather than later.
  - Everyday, there is a chance that tomorrow won't come.
  - The total discounted reward is finite.

Make Sure the policy is Finite

## -Type of MDP

A fully-observable MDP

The agent knows what state it is in.

A partially observable MDP (POMDP)

combines a MDP and a hidden Markov model.

The agent does not know what state it is in, but it can get some noisy signal of the state.

## -Grid World

Eg.

- There are four actions: up, down, left, and right.
- Every action is possible in every state.



The transition model  $P(s'|s, a)$

An action achieves its intended effect with probability 0.8.  
 An action leads to a 90-degree left turn with probability 0.1.  
 An action leads to a 90-degree right turn with probability 0.1.  
 If the robot bumps into a wall, it stays in the same square.

The reward function  $R(s)$  is the reward of entering state  $s$ .

$R(s_{01}) = -1$

$R(s_{02}) = 1$

Otherwise,  $R(s) = -0.04$

CQ: The robot is in  $s_{14}$  and tries to move to our right, what is the probability that the robot stays in  $s_{14}$ ?

- (A) 0.1  
 (B) 0.2  
 (C) 0.8  
 (D) 0.9  
 (E) 1.0

|   |       |   |   |    |
|---|-------|---|---|----|
| 1 | Start | 2 | 3 | 4  |
| 2 | X     |   |   |    |
| 3 |       |   |   | +1 |

Tips:

If environment is deterministic, don't worry about Agent would not land unexpected (100%)

## -A fixed Sequence of Action:

CQ: A fixed sequence of actions

- The robot could reach  $S_0$  w prob 0.15  
 OR The robot could reach  $S_0$  w prob  $0.1 \times 0.1 \times 0.8 \times 0.8 \times 0.8$

Q: Consider the action sequence "down, down, right, right, and right". This action sequence could take the robot to more than one square with positive probability.

- (A) True  
 (B) False  
 (C) I don't know

|   |       |   |   |    |
|---|-------|---|---|----|
| 1 | Start | 2 | 3 | 4  |
| 2 | X     |   |   | -1 |
| 3 |       |   |   | +1 |

So, fixed Sequence not fitting to MDP.

## -Policies

A policy specifies what the agent should do as a function of the current state.

A policy is

- non-stationary if it is a function of the state and the time.
- stationary if it is a function of the state.

## -The optimal policies - $R(s)$

The optimal policy of the grid world changes based on  $R(s)$  for any non-goal state  $s$ . It shows a careful balancing of risk and reward.

|   |       |      |    |    |
|---|-------|------|----|----|
| 1 | Start | 2    | 3  | 4  |
| 2 | X     | Risk | -1 | +1 |
| 3 |       |      |    |    |

- Optimal policy - Given  $V^*$

$$R(s) \rightarrow -0.04, \gamma = 1$$

The Expected Utility of a Policy

$$R(s) \text{ one-time reward of entering state } s.$$

$\rightarrow$   $V(s)$  expected utility of entering state  $s$  and following the policy thereafter.  $\downarrow$   
 the expected value of total discounted reward.

$\rightarrow$   $V(s)$  expected utility of entering state  $s$  and following the optimal policy  $\pi^*$  thereafter.

$\pi \rightarrow$  Given policy  
 $\pi^* \rightarrow$  Optimal Policy

$\rightarrow$  close to  $(+1)$ , higher expected utilities

|       |       |       |       |
|-------|-------|-------|-------|
| 1     | 2     | 3     | 4     |
| 0.705 | 0.655 | 0.611 | 0.338 |
| 0.762 | X     | 0.660 | -1    |
| 0.812 | 0.868 | 0.918 | +1    |

Figure:  $V^*(s)$  for  $\gamma = 1$  and  $R(s) = -0.04$ , vs  $s \neq s_3$ .

$\rightarrow$  Quantity  $Q^*$ :

Expected Utilities of taking actions in current state

$$Q^*(s_a) = \sum_{s'} P(s'|s, a) V^*(s')$$

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

Eg.

$$\begin{aligned} Q(S_0, down) &= 0.8 \times 0.655 + 0.1 \times 0.338 + 0.1 \times 0.338 = 0.7525 \\ Q(S_0, left) &= 0.8 \times 0.660 + 0.1 \times 0.611 + 0.1 \times 0.611 = 0.6571 \\ Q(S_0, right) &= 0.8 \times 0.338 + 0.1 \times 0.611 + 0.1 \times 0.611 = 0.4375 \\ Q(S_0, up) &= 0.8 \times 0.611 + 0.1 \times 0.655 + 0.1 \times 0.338 = 0.5931 \\ T(S_0) &= \text{left} \end{aligned}$$

|       |       |       |       |
|-------|-------|-------|-------|
| 1     | 2     | 3     | 4     |
| 0.705 | 0.655 | 0.611 | 0.338 |
| 0.762 | X     | 0.660 | -1    |

## -Solving $V^*(s)$ Using Value Iteration.

$V$ : EU of Agent  $\rightarrow$   $\pi$

then follow policy after get  $\pi$

immediate reward (long time Total discount reward)

Q: long time total discount reward, but we are already in a state what happen if we take particular action.

When  $R(s) > 0$ , what does the optimal policy look like?

|   |       |   |   |   |
|---|-------|---|---|---|
| 1 | Start | 2 | 3 | 4 |
| 2 | X     |   |   |   |
| 3 |       |   |   |   |

- avoids both goal states.

# X. V & Q relation (Bellmen)

V and Q are defined recursively in terms of each other.

$$\textcircled{Q} \quad V^*(s) = R(s) + \gamma \max_a Q^*(s, a) \quad (3)$$

$$\textcircled{Q'} \quad Q^*(s, a) = \sum_{s'} P(s'|s, a) V^*(s') \quad (4)$$

Combining equations 3 and 4, we get the Bellman equations:

Bellmen:

$$V^*(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^*(s') \quad (5)$$

$V^*(s)$  are the unique solution to the Bellman equations.

Eg.

Write down  $V^*(s_{11})$

|   | 1     | 2     | 3     | 4     |
|---|-------|-------|-------|-------|
| 1 | 0.705 | 0.655 | 0.611 | 0.388 |
| 2 | 0.762 | X     | 0.660 | -1    |
| 3 | 0.812 | 0.868 | 0.918 | +1    |

Write down the Bellman equation for  $V^*(s_{11})$ .

$$V(s_{11}) = -0.04 + \gamma \max_a [0.8 \cdot V(s_{12}) + 0.1 \cdot V(s_{21}) + 0.1 \cdot V(s_{31})], \text{ (right)} \\ 0.9 \cdot V(s_{11}) + 0.1 \cdot V(s_{12}), \text{ (up)} \\ 0.9 \cdot V(s_{11}) + 0.1 \cdot V(s_{21}) + 0.1 \cdot V(s_{31}) \text{ (down)}$$

## Value Iteration Algo

Tips: Solve System of Bellman equation efficiently? No, max  
So, Iteration needed

### Solving for $V^*(s)$ iteratively

The Bellman equations:

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

Let  $V_i(s)$  be the values for the i-th iteration.

1. Start with arbitrary initial values for  $V_0(s)$ .
2. At the i-th iteration, compute  $V_{i+1}(s)$  as follows.

$$\rightarrow V_{i+1}(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V_i(s')$$

3. Terminate when  $\max_s |V_i(s) - V_{i+1}(s)|$  is small enough.

If we apply the Bellman update infinitely often, the  $V_i$ 's are guaranteed to converge to the optimal values.

Eg.

CQ: Calculating  $V_1(s_{23})$

$$V(s_{23}) = -0.04 + \gamma \max_a [0.8 \cdot 0 + 0.1 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0], \text{ right} \\ 0.8 \cdot 0 + 0.1 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0, \text{ up} \\ 0.8 \cdot 0 + 0.1 \cdot 0 + 0.1 \cdot 0 + 0.1 \cdot 0, \text{ down}$$

CQ: What is  $V_1(s_{23}) = -0.04$

- (A)  $(-\infty, 0]$    (B)  $[0, 0.25]$    (C)  $[0.25, 0.5]$   
(D)  $[0.5, 0.75]$    (E)  $[0.75, 1]$

$V_0(s)$ :

|   | 1 | 2 | 3 | 4  |
|---|---|---|---|----|
| 1 | 0 | 0 | 0 | 0  |
| 2 | 0 | X | 0 | -1 |
| 3 | 0 | 0 | 0 | +1 |

CQ: Calculating  $V_1(s_{33})$

$$V_1(s_{33}) = -0.04 + \gamma \max_a [0.8 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0], \text{ right} \\ 0.8 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0, \text{ up} \\ 0.8 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 0, \text{ down}$$

$$= -0.04 + 0.8 = 0.76$$

$\textcircled{Step}$

CQ: What is  $V_1(s_{33})$ ?

- (A) 0.26   (B) 0.36   (C) 0.46  
(D) 0.56   (E) 0.76

$V_0(s)$ :

|   | 1 | 2 | 3 | 4  |
|---|---|---|---|----|
| 1 | 0 | 0 | 0 | 0  |
| 2 | 0 | X | 0 | -1 |
| 3 | 0 | 0 | 0 | +1 |

$V_1(s)$ :

|   | 1     | 2     | 3     | 4     |
|---|-------|-------|-------|-------|
| 1 | -0.04 | -0.04 | -0.04 | -0.04 |
| 2 | -0.08 | X     | -0.04 | -1    |
| 3 | -0.08 | -0.04 | 0.76  | +1    |

$V_2(s)$ :

|   | 1     | 2     | 3     | 4     |
|---|-------|-------|-------|-------|
| 1 | -0.08 | -0.08 | -0.08 | -0.08 |
| 2 | -0.08 | X     | 0.464 | -1    |
| 3 | -0.08 | 0.56  | 0.832 | +1    |

Thus, Each state accumulates negative rewards until algo finds a path to +1

synchronous:  $\boxed{\text{II}} \rightarrow \boxed{\text{II}}$  — update  $\xrightarrow{\text{same time}}$   
asynchronous: real time update,  
converges faster

## Policy Iteration Algo (Reforcement needed)

Inspired: to derive the optimed  
Policy doesn't require accurate  
estimate of utility function  $V^*(s)$

### Policy Iteration we got

► Policy iteration alternates between two steps.

- Policy evaluation: Given a policy  $\pi_i$ , calculate  $V^{\pi_i}(s)$ , which is the utility of each state if  $\pi_i$  were to be executed.

- Policy improvement: Calculate a new policy  $\pi_{i+1}$  using  $V^{\pi_i}$ .

Terminates when there is no change in the policy.

policy iteration:  $\pi_i(s) \rightarrow V^{\pi_i}(s) \rightarrow \pi_b(s) \rightarrow V^{\pi_b}(s) \rightarrow \dots \rightarrow \pi^*(s) \rightarrow V^{\pi^*}(s) \rightarrow \pi^{**}(s)$

(Policy  $\boxed{\text{II}} \rightarrow \Delta \rightarrow \Delta \rightarrow \Delta$ ) (Value  $\boxed{\Delta} \rightarrow \Delta$ )

► Policy improvement:  $V^{\pi_i}(s) \rightarrow \pi_{i+1}(s)$

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} P(s'|s, a) V^{\pi_i}(s').$$

► Policy evaluation:  $\pi_i(s) \rightarrow V^{\pi_i}(s)$

$$V_i(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) V_i(s')$$

Simplified Bellmen without  
max

Policy evaluation:  $V(S_{ii}) = -0.04 + \gamma (0.8 V(S_{ii}) + 0.1 V(S_{ii}) + 0.1 V(S_{ii}))$

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

$$V_{ii} = -0.04 + \gamma (0.8 \cdot V_{ii} + 0.1 \cdot V_{ii} + 0.1 \cdot V_{ii})$$

$$0.9 \cdot V_{ii} + 0.1 \cdot V_{ii} + 0.1 \cdot V_{ii}$$

$$0.8 \cdot V_{ii} + 0.2 \cdot V_{ii} + 0.1 \cdot V_{ii}$$

Bellman equations:

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

Eg.

$$\textcircled{0} \quad \begin{array}{|c|c|} \hline 1 & 2 \\ \hline \boxed{1} & +1 \\ \hline 2 & -1 \\ \hline \end{array}$$

Policy  $\pi(S_{ii}) = \text{right}$   $\pi(S_{ii}) = \text{right}$

Policy evaluation

$$\left\{ \begin{array}{l} V(S_{ii}) = -0.04 + 0.8 (+) + 0.1 V(S_{ii}) + 0.1 V(S_{ii}) \\ V(S_{ii}) = -0.04 + 0.8 \cdot (-) + 0.1 V(S_{ii}) + 0.1 V(S_{ii}) \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} 0.9 V(S_{ii}) - 0.1 V(S_{ii}) = 0.76 \\ -0.1 V(S_{ii}) + 0.9 V(S_{ii}) = -0.84 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} 8 V(S_{ii}) = 6 \Rightarrow V(S_{ii}) = 0.75 \\ 0.9 V(S_{ii}) = -0.765 \Rightarrow V(S_{ii}) = -0.85 \end{array} \right.$$

$\textcircled{2}$

Policy

Improvement

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline \boxed{1} & +1 \\ \hline 2 & -1 \\ \hline \end{array}$$

$\pi(S_{ii}) = ?$

$$\rightarrow \begin{array}{l} \text{right: } 0.8 (+) + 0.1 (0.75) + 0.1 (-0.85) = -0.85 \\ 0.8 \quad 0.075 \quad -0.085 \\ \text{up: } 0.8 (0.75) + 0.1 (-0.85) + 0.1 (+) = 0.775 \\ 0.6 \quad 0.075 \quad 0.1 \\ \text{left: } 0.8 (0.75) + 0.1 (-0.85) + 0.1 (0.75) = 0.59 \\ 0.6 \quad -0.085 \quad 0.075 \\ \text{down: } 0.8 (-0.85) + 0.1 (+) + 0.1 (0.75) = -0.505 \\ -0.68 \quad 0.1 \quad 0.075 \end{array}$$

$\pi(S_{ii}) = \text{right}$

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline \boxed{1} & +1 \\ \hline 2 & -1 \\ \hline \end{array}$$

$\pi(S_{ii}) = ?$

$$\begin{array}{l} \text{right: } 0.8 (-) + 0.1 (0.75) + 0.1 (-0.85) = -0.85 \\ 0.8 \quad 0.075 \quad -0.085 \\ \text{up: } 0.8 (0.75) + 0.1 (-0.85) + 0.1 (-) = -0.185 \\ 0.6 \quad -0.085 \quad -0.1 \\ \text{left: } 0.8 (-0.85) + 0.1 (-0.85) + 0.1 (0.75) = -0.89 \\ -0.68 \quad -0.085 \quad 0.075 \\ \text{down: } 0.8 (-0.85) + 0.1 (-) + 0.1 (-0.85) = -0.865 \\ -0.68 \quad -0.1 \quad -0.085 \end{array}$$

$\pi(S_{ii}) = \text{up}$

$\textcircled{3}$

previous policy

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline \boxed{1} & +1 \\ \hline 2 & -1 \\ \hline \end{array}$$

new policy

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline \boxed{1} & +1 \\ \hline 2 & \uparrow -1 \\ \hline \end{array}$$

Policy  $\pi(S_{ii}) = \text{right}$   $\pi(S_{ii}) = \text{up}$

Policy evaluation

$$\left\{ \begin{array}{l} V(S_{ii}) = -0.04 + 0.8 (+) + 0.1 V(S_{ii}) + 0.1 V(S_{ii}) \\ V(S_{ii}) = -0.04 + 0.8 V(S_{ii}) + 0.1 V(S_{ii}) + 0.1 (-) \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} 0.9 V(S_{ii}) + 0.1 V(S_{ii}) = 0.76 \\ -0.1 V(S_{ii}) + 0.9 V(S_{ii}) = -0.14 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} 7.3 V(S_{ii}) = 6.7 \Rightarrow V(S_{ii}) = 0.918 \\ 0.9 V(S_{ii}) = 0.5944 \Rightarrow V(S_{ii}) = 0.660 \end{array} \right.$$

Terminate until no change

# Reinforcement Learning

## Eg.

### An Reinforcement Learning Agent

Let's consider fully observable, single-agent reinforcement learning. We will formalize this problem as a Markov decision process.

- Given the possible states and the set of actions.
- Observes the state and the rewards received.
- Carries out an action.
- Goal is to maximize its discounted reward.

Why is reinforcement learning challenging?

- Which action was responsible for this reward/punishment?
- How will this action impact my utility?
- Should I explore or exploit?

## Model-based

### Include $P(s'|s,a)$ and $R(s)$

#### Passive Learning Agent

- Fix a policy  $\pi$ .
- Goal is to learn  $V^\pi(s)$  (the expected value of policy  $\pi$  for state  $s$ ).
- Similar to policy evaluation.
- Does not know the transition model  $P(s'|s,a)$  nor the reward function  $R(s)$ .
- Solution: Adaptive Dynamic Programming
  - Learn  $P(s'|s,a)$  and  $R(s)$  using the observed transitions and rewards.
  - Learn  $V^\pi(s)$  by solving Bellman equations (exactly or iteratively).

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

## Algo step:

#### Passive ADP Algorithm

- Repeat steps 2 to 5.
- Follow policy  $\pi$  and generate an experience  $(s, a, s', r')$ .
- Update reward function:  $R(s') \leftarrow r'$ .
- Update the transition probability.

$$N(s, a) = N(s, a) + 1$$

$$N(s, a, s') = N(s, a, s') + 1$$

$$P(s'|s, a) = N(s, a, s') / N(s, a)$$

- Derive  $V^\pi(s)$  by using the Bellman equations.

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

Eg.

|          |    |
|----------|----|
| $s_{11}$ | +1 |
| $s_{21}$ | -1 |

- $\pi(s_{11}) = \text{down}, \pi(s_{21}) = \text{right}$
- $\gamma = 0.9$
- $R(s_{11}) = -0.01, R(s_{21}) = -0.01, R(s_{12}) = 1, R(s_{22}) = -1$
- $N(s, a) = 5, \forall s, a$
- $N(s, a, s') = 3$  for the intended direction.
- $N(s, a, s') = 1$  for a direction to the left or right of the intended direction.
- The current state is  $s_{11}$ .

experience

|                                  |    |
|----------------------------------|----|
| $s_{11}$ , down, $S_{11}, -0.01$ | +1 |
| $s_{21}$                         | -1 |

- No need to update the reward function.

- Update the counts.

$$N(s_{11}, \text{down}) = 6 \text{ and } N(s_{11}, \text{down}, s_{21}) = 4.$$

- Solve the Bellman equations.  $P(s_{11} | S_{11}, \text{down}) = \frac{4}{6} = 0.667$ .

$$V(s_{11}) = -0.04 + 0.9(0.667V(s_{21}) + 0.167(-1) + 0.167V(s_{12}))$$

$$V(s_{21}) = -0.04 + 0.9(0.6(-1) + 0.2V(s_{11}) + 0.2V(s_{22}))$$

The solutions are:

$$V(s_{11}) = -0.4378, V(s_{21}) = -0.8034$$

## → Active ADP

short term: Reward from action

long term: get data learn model

The passive ADP agent learns the expected value of a fixed policy.

What action should the agent take at each step?

Two things are useful for the agent to do:

① exploit: take an action that maximizes  $V(s)$ .

② explore: take an action that is different from the optimal one.

The greedy agent seldom converges to the optimal policy and sometimes converges to horrible policies because the learned model is not the same as the true environment.

There is a trade-off between exploitation and exploration.

#### Trade off Exploitation and Exploration

①  $\epsilon$ -greedy exploration strategy:

- Select random action with probability  $\epsilon$  and
- Select the best action with probability  $1 - \epsilon$ .
- may decrease  $\epsilon$  over time.

② Softmax selection using Gibbs/Boltzmann distribution:

- Choose action  $a$  with probability  $\frac{Q(s, a) / T}{\sum_a Q(s, a) / T}$ .
- $T > 0$  is the temperature. When  $T$  is high, the distribution is close to uniform. When  $T$  is low, the higher-valued actions have higher probabilities.

③ Initialize the values optimistically to encourage exploration.



#### Optimistic Utility Values to Encourage Exploration

We will learn  $V^+(s)$  (the optimistic estimates of  $V(s)$ ).

$$V^+(s) \leftarrow R(s) + \gamma \max_a f\left(\sum_{s'} P(s'|s, a) V^+(s'), N(s, a)\right)$$

$$f(u, n) = \begin{cases} R(s), & \text{if } n < N_e \\ u, & \text{otherwise} \end{cases}$$

$f(u, n)$  trade-offs exploitation and exploration.

④  $G^+$  is the optimistic estimate of the best possible reward obtainable in any state.

⑤ If we haven't visited  $(s, a)$  at least  $N_e$  times, assume its expected value is  $R^+$ .

Otherwise, use the current  $V^+(s)$  value.

#### Active ADP Algorithm

- Initialize  $R(s), V^+(s), N(s, a), N(s, a, s')$ .

- Repeat steps 3 to 7 until we have visited each  $(s, a)$  check at least  $N_e$  times and the  $V^+(s)$  values converged.

- Determine the best action  $a$  for current state  $s$  using  $V^+(s)$ .

$$a = \arg \max_a f\left(\sum_{s'} P(s'|s, a) V^+(s'), N(s, a)\right), f(u, n) = \begin{cases} R^+, & \text{if } n < N_e \\ u, & \text{otherwise} \end{cases}$$

- Take action  $a$  and generate an experience  $(s, a, s', r')$ .

- Update reward function:  $R(s') \leftarrow r'$

- Update the transition probability.

$$N(s, a) = N(s, a) + 1, N(s, a, s') = N(s, a, s') + 1$$

$$P(s'|s, a) = N(s, a, s') / N(s, a)$$

- Updated  $V^+(s)$  using the Bellman updates. check convergence.

$$V^+(s) \leftarrow R(s) + \gamma \max_a f\left(\sum_{s'} P(s'|s, a) V^+(s'), N(s, a)\right)$$

ADP: more efficient experience, fast learn

Q: low cost, computation

## → Model-free.

## → Temporal Difference

Introduce it:

#### Bellman Equations for $Q(s, a)$

$Q(s, a)$  is the expected value of performing action  $a$  in state  $s$ . We can define Bellman equations for both  $V(s)$  and  $Q(s, a)$ .

#### Bellman equations for $V(s)$ :

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

#### Bellman equations for $Q(s, a)$ :

$$\rightarrow Q(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$$

Learning  $V(s)$  and  $Q(s, a)$  are equivalent! What is the advantage of learning  $Q(s, a)$ ?

#### Temporal Difference Error

Assume that we observed  $(s_1, r_1, a_1, s_2)$ . Based on this transition, what should  $Q(s_1, a_1)$  satisfy?

Start with the Bellman equations for  $Q(s_1, a_1)$ .

$$Q(s_1, a) = R(s_1) + \gamma \sum_{s'} P(s'|s_1, a_1) \max_{a'} Q(s', a')$$

$Q(s_1, a)$  should be computed by the RHS of the above equation. Assume that this transition always occurs  $(P(s_2, s_3, a_2) = 1)$ . Thus,  $Q(s_1, a)$  should be

$$R(s_1, a) + \max_{a'} Q(s_2, a')$$

#### Temporal difference (TD) error:

$$(R(s_1) + \gamma \max_{a'} Q(s_2, a')) - Q(s_1, a)$$

1. 从最近开始，移到第一个种子。  
2. 一旦遇到和前一个相同的种子，自动生成一个即时奖励，并且根据到目前为止的值。  
3. 价值更新时将新的价值加到前一个种子的值上。

价值更新规则：

$$V(s) := V(s) + \alpha(r + \gamma V(s') - V(s))$$

→ Q learning from TD Updates

#### Q-Learning Updates

TD

Given an experience  $(s, r, a, s')$ , update  $Q(s, a)$  as follows

$$Q(s, a) \leftarrow Q(s, a) + \alpha (R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

predicted value  
current value

An alternative version:

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + \alpha (R(s) + \gamma \max_{a'} Q(s', a'))$$

current value  
predicted value

#### Passive Q-Learning Algorithm

- Repeat steps 2 to 4.
- Follow policy  $\pi$  and generate an experience  $(s, r, a, s')$ .
- Update reward function:  $R(s) \leftarrow r$ .
- Update  $Q(s, a)$  by using the temporal difference update rules:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

The learning rate  $\alpha$ :  $\alpha$  controls the size of each update. If  $\alpha$  decreases as  $N(s, a)$  increases,  $Q$  values will converge to the optimal values. For example,  $\alpha(N(s, a)) = \frac{1}{\sqrt{N(s, a)}}$ .

#### Active Q-Learning Algorithm

→ QL, no policy but

- Initialize  $R(s), Q(s, a), N(s, a), N(s, a, s')$ .
- Repeat steps 3 to 6 until we have visited each  $(s, a)$  at least  $N_e$  times and the  $Q(s, a)$  values converged.
- Determine the best action  $a$  for current state  $s$  using  $V^+(s)$ .
- $a = \arg \max_a f\left(\sum_{s'} P(s'|s, a) V^+(s'), N(s, a)\right), f(u, n) = \begin{cases} R^+, & \text{if } n < N_e \\ u, & \text{otherwise} \end{cases}$
- Take action  $a$  and generate an experience  $(s, r, a, s')$ .
- Update reward function:  $R(s) \leftarrow r$ .
- Update  $Q(s, a)$  using the temporal difference update rules.

$$Q(s, a) \leftarrow Q(s, a) + \alpha (R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$



## Independence

## - Unconditional Independence

Definition ((unconditional) independence)

$X$  and  $Y$  are (unconditionally) independent iff

$$P(X|Y) = P(X)$$

$$P(Y|X) = P(Y)$$

$$P(X \wedge Y) = P(X)P(Y)$$

Learning  $Y$  does NOT influence your belief about  $X$ .

## - Conditional Independence

Definition (conditional independence)

$X$  and  $Y$  are conditionally independent given  $Z$  if

$$P(X|Y|Z) = P(X|Z).$$

$$P(Y|X|Z) = P(Y|Z).$$

$$P(Y \wedge X|Z) = P(Y|Z)P(X|Z).$$

Learning  $Y$  does NOT influence your belief about  $X$  if you already know  $Z$ .

Eg.

CQ: Deriving a compact representation

of random variables,  $2^n - 1$  probabilities.

Boolean

CQ: Consider a model with three random variables,  $A, B, C$ .  
 1. What is the minimum number of probabilities required to specify the joint distribution?  $7$  probabilities.  
 $P(A \wedge B \wedge C) = P(A)P(B|A)P(C|A \wedge B)$   
 $= P(A)P(B|A)P(C|A \wedge B)$   
 2. Assume that  $A, B$ , and  $C$  are independent. What is the minimum number of probabilities required to specify the joint distribution?  $3$  probabilities.  
 $P(A \wedge B \wedge C) = P(A)P(B)P(C)$   
 $= P(A)P(B|A)P(C|A \wedge B)$   
 $= P(A)P(B|A)P(C|B)$   
 $= P(A)P(B|A)P(C|B)$   
 $\text{A} \quad \text{B} \quad 1+1+1=3 \quad 1+2+4=7$

3. Assume that  $A$  and  $B$  are conditionally independent given  $C$ . What is the minimum number of probabilities required to specify the joint distribution?  $5$  probabilities

$$\begin{aligned} P(A \wedge B \wedge C) &= P(C)P(A|C)P(B|CNA) \\ &= P(C)P(A|C)P(B|CNA) \\ &= P(C) + P(A|C)P(B|C) \\ &= P(C) + P(A|C)P(B|CNA) \\ &= P(C) + P(A|C)P(B|CNA) \\ &= 1+2+2=5 \end{aligned}$$

## Bayesian Network

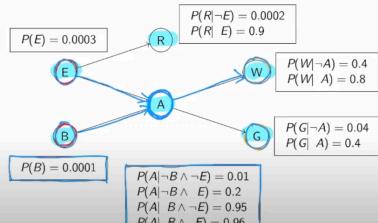
A Bayesian Network

- is a compact version of the joint distribution, and
- takes advantage of the (unconditional and conditional) independence among the variables.

Eg. no circle,  $\Rightarrow$

A Bayesian Network for the Holmes Scenario

$$2^6 - 1 = 63 \quad 1+1+2+4+2+0 = 12$$



## - Representing joint distribution

Can Compute each joint  $P$  use formula:

$$P(X_1 \wedge \dots \wedge X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

Eg.

Example: What is the probability that

- The alarm has sounded,
- Neither a burglary nor an earthquake has occurred,
- Both Watson and Gibbon call and say they hear the alarm, and
- There is no radio report of an earthquake?

$$\begin{aligned} P(\neg B \wedge \neg E \wedge \neg R \wedge \neg A \wedge \neg G) \\ = P(\neg B)P(\neg E)P(\neg A \wedge \neg R \wedge \neg E)P(\neg R \wedge \neg E)P(\neg W|A)P(\neg G|A) \\ = (-1.0001)(-1.0003)(0.01)(-1.0002)(0.4)(0.8) \\ = 0.0032 \end{aligned}$$

## - Un/Conditional Independence relationship in Bayes

Eg. Causal chain



Burglary, Watson independent?

No

B, W Conditionally inde?

Yes

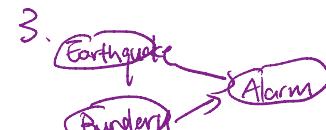


Watson, Gibbon (uncor) independent?

No

Watson, Gibbon can inde?

Yes



Earthquake, Burglary independent?

Yes

CQ: Are Earthquake and Burglary conditionally independent given Alarm?

the "explaining away" effect



Suppose that the Alarm is going off.

- (A) Yes  
 (B) No

If an earthquake is happening, then it is less likely that the Alarm is caused by a Burglary.

If an earthquake is not happening, then it is more likely that a burglary is happening and causing the alarm to go off.

## - Bayes Net : D-Separation

(Independent or Not)

Definition (D-Separation)

$E$  d-separates  $X$  and  $Y$   
 iff  $E$  blocks every un-directed path between  $X$  and  $Y$ .

If  $E$  d-separates  $X$  and  $Y$ ,

then  $X$  and  $Y$  are independent given  $E$ .

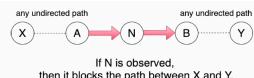
path we consider

① Un-direct

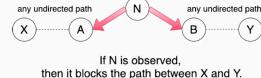
② multiple paths

③ multiple node

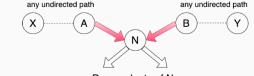
$X \dashv \vdash [S e r v i c e] \dashv \vdash Y$



If  $N$  is observed, then it blocks the path between  $X$  and  $Y$ .



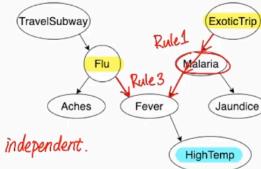
If  $N$  is observed, then it blocks the path between  $X$  and  $Y$ .



If  $N$  and  $N$ 's descendants are NOT observed, then they block the path between  $X$  and  $Y$ .

Eg.

CQ 3b: Are Flu and ExoticTrip conditionally independent given HighTemp?



Not independent.

ExoticTrip -> Malaria -> Fever -> Jaundice -> HighTemp

# - Construction Bayes Net

For a joint probability distribution, there are multiple correct Bayesian networks.

A Bayesian network is correct if every independence relationship the network represents is correct.  $\rightsquigarrow$  exists in the joint distribution.

We prefer one Bayesian network over another one if the former requires fewer probabilities.

Joint P dist  $\{BN_1, BN_2, BN_3\}$ ,  $BN_1$  with fewer P is good,

Inde correct  $\rightarrow$  Network correct

1. Determine the set of variables for the domain.

2. Order the variables  $(X_1, \dots, X_n)$ .

3. For each variable  $X_i$  in the ordering,

3.1 Choose the node's parents:

Choose the smallest set of parents from  $\{X_1, \dots, X_{i-1}\}$  such that given  $Parents(X_i) = X_j$ , independently of all the nodes in  $\{X_1, \dots, X_{i-1}\} - Parents(X_i)$ . Formally,

$$P(X_i | Parents(X_i)) = P(X_i | X_{i-1} \wedge \dots \wedge X_1)$$

3.2 Create a link from each parent of  $X_i$  to the node  $X_i$ .

3.3 Write down the conditional probability table  $P(X_i | Parents(X_i))$ .

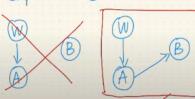
Eg.

Consider the Bayesian network.



Construct a correct Bayesian network based on the variable ordering: W, A, and B.

Step 1: add (W) to the network.

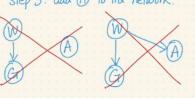


Consider the Bayesian network:



Construct a correct Bayesian network based on the variable ordering: W, G, and A.

Step 1: add (W) to the network.



Step 1: (W)  
Step 2: add (A) to the network.  
~~(W)  $\rightarrow$  (A)~~  
 $(W) \rightarrow (A)$

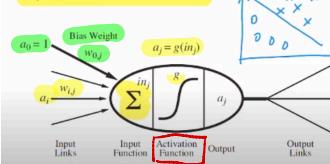
Step 1 : (W)  
Step 2 : add (G) to the network.  
~~(W)  $\rightarrow$  (G)~~  
 $(W) \rightarrow (G)$

Step 1 : (W)  
Step 2 : add (G) to the network.  
~~(W)  $\rightarrow$  (G)~~  
 $(W) \rightarrow (G)$   
~~(W)  $\rightarrow$  (A)~~  
 $(W) \rightarrow (A)$

make it more Specific

# 0 Neural Network - mathematical model

A linear classifier  $\rightarrow$  it "fires" when a linear combination of its inputs exceeds some threshold.



Neuron  $j$  computes a weighted sum of its input signals.

$$i_j = \sum_{i=0}^n w_{ij} a_i$$

Neuron  $j$  applies an activation function  $g$  to the weighted sum to derive the output.  $a_j = g(i_j) = g(\sum_{i=0}^n w_{ij} a_i)$ .

## About Activation Function:

What are some desirable properties of the activation function?

- It should be non-linear:  
combining linear functions will not produce a non-linear function.  
complex relationships are often non-linear.
- It should mimic the behaviour of real neurons.  
If the weighted sum of input signals is large enough, the neuron fires. Otherwise, it does not fire.
- It should be differentiable almost everywhere.  
learn a neural network using gradient descent, which requires the activation function to be differentiable.

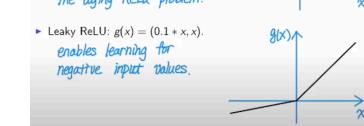
## Common Activation Function:

Step function:  $g(x) = 1$  if  $x > 0$ .  $g(x) = 0$  if  $x \leq 0$ .  
simple to use, not differentiable  
not used in practice, but useful  
to explain concepts.

Sigmoid function:  $g(x) = \frac{1}{1 + e^{-x}}$   
can approximate the step function.  
clear, bounded prediction.  
differentiable.  
"vanishing gradient" problem.

Rectified linear unit (ReLU):  $g(x) = \max(0, x)$ .  
computationally efficient.  
non-linear and differentiable.  
the dying ReLU problem.

Leaky ReLU:  $g(x) = (0.1 + x) \cdot x$ .  
enables learning for negative input values.



## network

Connecting the neurons together into a network

Feed-forward network  
directed acyclic graph  
no loops.  
a function of its inputs.

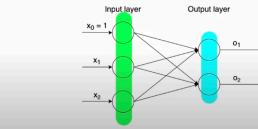
Recurrent network  
has loops  
has memory.  
better model of human brain,  
but more difficult to interpret and train.

# - Perceptrons

Single-layer feed-forward neural network

The inputs are connected directly to the outputs.

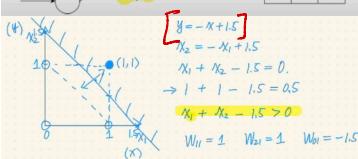
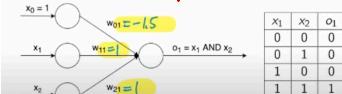
Can represent logical functions, e.g. AND, OR, and NOT.



## Eg. logical GATE

CQ: Consider the perceptron below where the activation function is the step function ( $g(x) = 1$  if  $x > 0$ ,  $g(x) = 0$  if  $x \leq 0$ ). What should the weights  $w_{01}, w_{11}$  and  $w_{21}$  be such that the perceptron represents an AND function?

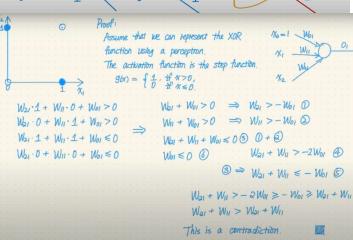
$$w_{01} + w_{11} + w_{21} \dots$$



## Eg.

CQ: Why can't a perceptron represent XOR?

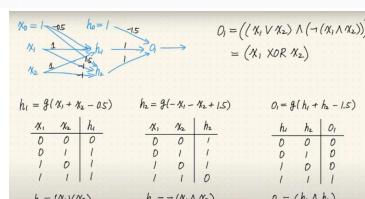
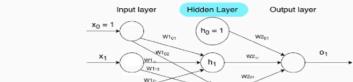
① a perceptron is a linear classifier.  
XOR is not linearly separable.



XOR as a 3-Layer Neural Network

$$\textcircled{2} \quad O_1 = ((X_1 \vee X_2) \wedge (\neg(X_1 \wedge X_2)))$$

Can you come up with the weights such that the following network represents the XOR function?



# Gradient Descent

## Gradient Descent

expected outputs based on training data  
actual outputs based on training data  
 $\hat{f}(x)$  network.

"Walking downhill and always taking a step in the direction that goes down the most."

error

A local search algorithm to find the minimum of a function.

Steps of the algorithm:

- Initialize weights randomly.
- Change each weight in proportion to the negative of the partial derivative of the error with respect to the weight.

$$w := w - \eta \frac{\partial \text{error}}{\partial w}$$

- $\eta$  is the learning rate
- Sum over all training examples.

Terminate after some number of steps

when the error is small or when the changes get small.

## Weight

How do we update the weights based on the data points?  
cheaper steps, weights become more accurate more quickly.  
not guaranteed to converge.

Gradient descent updates the weights after sweeping through all the examples.

To speed up learning, update weights after each example.

Incremental gradient descent: update weights after each example.

Stochastic gradient descent: same as incremental gradient descent except that each example is chosen randomly.

Trade off learning speed and convergence.

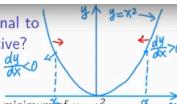
Batched gradient descent:

update weights after a batch of examples.

start w/ small batches to learn quickly.  
then, increase batch size so it converges.

## Tips: How to derive Gradient Descent

Why update the weight proportional to the negative of the partial derivative?



Suppose that we want to find the minimum of  $y = x^2$ .

Start with  $x = x_0$ .

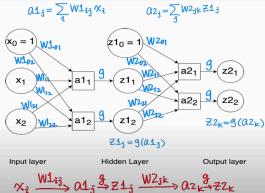
In what direction should we change the value of  $x$ ?  
change  $x$  in the direction of the negative of the partial derivative.

By what amount should we change the value of  $x$ ?  
What is the step size?

Take a step in proportion to the magnitude of the partial derivative.

## The Back propagation Algorithm

### A 2-Layer Neural Network



## The Back-propagation Algorithm

An efficient method of calculating the gradients in a multi-layer neural network.

There are some training examples  $(\bar{x}_i, \bar{y}_i)$  and an error/loss function  $E(z_i, \bar{y}_i)$ . Perform 2 passes.

Forward pass: compute the error  $E$  given the inputs and the weights.

Backward pass: compute the gradients  $\frac{\partial E}{\partial W_{2jk}}$  and  $\frac{\partial E}{\partial W_{1ij}}$ .

Update each weight by the sum of the partial derivatives for all the training examples.

## Forward Pass

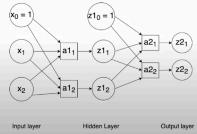
### Forward Pass for a 2-layer Network

Calculate the values of  $z_{1j}$  and  $z_{2k}$  and  $E$ .

$$a_{1j} = \sum_i x_i W_{1ij} \quad z_{1j} = g(a_{1j}) \quad (1)$$

$$a_{2k} = \sum_j z_{1j} W_{2jk} \quad z_{2k} = g(a_{2k}) \quad (2)$$

$$E(z_2, y) \quad (3)$$



Input layer      Hidden Layer      Output layer

## Backward

### Backward Pass for a 2-layer Network

Calculate the gradients for  $W_{1ij}$  and  $W_{2jk}$ .

$$\frac{\partial E}{\partial W_{2jk}} = \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{1j}} \cdot \frac{\partial z_{1j}}{\partial a_{1j}} \cdot \frac{\partial a_{1j}}{\partial x_i} \quad (4)$$

$$\frac{\partial E}{\partial W_{1ij}} = \frac{\partial E}{\partial a_{1j}} \cdot \frac{\partial a_{1j}}{\partial x_i} \quad (5)$$



## Recursive relationship

For unit  $j$ ,  $\delta_j = \frac{\partial E}{\partial a_j}$

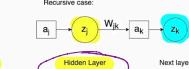
$$\delta_j = \begin{cases} \frac{\partial E}{\partial z_j} \times g'(a_j), & \text{base case: } j \text{ is an output unit,} \\ \left( \sum_k \frac{\partial E}{\partial W_{jk}} \right) \times g'(a_j), & \text{recursive case: } j \text{ is a hidden unit} \end{cases} \quad (6)$$

Base case:



Output layer

Recursive case:



Hidden Layer

Next layer

### Deriving the expressions for $\frac{\partial E}{\partial W_{2jk}}$

$$\frac{\partial E}{\partial W_{2jk}} = \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{1j}} \cdot \frac{\partial z_{1j}}{\partial a_{1j}} \cdot \frac{\partial a_{1j}}{\partial x_i}$$

$$\text{an example w/ } i=1 \text{ and } j=2$$

$$\frac{\partial E}{\partial W_{2jk}} = \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{12}} \cdot \frac{\partial z_{12}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial x_1}$$

chain rule  $\frac{d}{dx} f(f(x)) = \frac{df}{dx} \frac{df}{dx}$

$$E = \sum_k \frac{(z_{2k} - y_k)^2}{2} \Rightarrow \frac{\partial E}{\partial z_{2k}} = (z_{2k} - y_k)$$

$$g(x) = \frac{1}{1 + e^{-x}} \Rightarrow g'(x) = g(x)(1 - g(x))$$

$$a_{2k} = \sum_j W_{2kj} z_{1j} \Rightarrow \frac{\partial a_{2k}}{\partial z_{1j}} = W_{2kj}$$

$$z_{1j} = g(a_{1j}) \Rightarrow \frac{\partial z_{1j}}{\partial a_{1j}} = g'(a_{1j})$$

$$a_{1j} = \sum_i x_i W_{1ij} \Rightarrow \frac{\partial a_{1j}}{\partial x_i} = W_{1ij}$$

$$\frac{\partial E}{\partial W_{2jk}} = \sum_i \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{1j}} \cdot \frac{\partial z_{1j}}{\partial a_{1j}} \cdot \frac{\partial a_{1j}}{\partial x_i}$$

$$= \sum_i \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{1j}} \cdot \frac{\partial z_{1j}}{\partial a_{1j}} \cdot W_{1ij}$$

$$= \sum_i \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{1j}} \cdot \frac{\partial z_{1j}}{\partial a_{1j}} \cdot g'(a_{1j}) \cdot W_{1ij}$$

$$\text{an example w/ } i=1 \text{ and } j=2$$

$$\frac{\partial E}{\partial W_{2jk}} = \sum_k \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{12}} \cdot \frac{\partial z_{12}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial x_1}$$

$$= \sum_k \frac{\partial E}{\partial z_{2k}} \cdot \frac{\partial z_{2k}}{\partial a_{2k}} \cdot \frac{\partial a_{2k}}{\partial z_{12}} \cdot \frac{\partial z_{12}}{\partial a_{12}} \cdot g'(a_{12}) \cdot W_{12j}$$

$$\text{Defining the delta values}$$

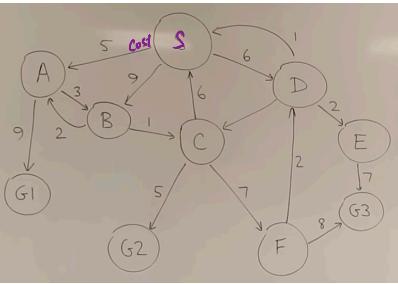
$$\frac{\partial E}{\partial W_{2jk}} = \frac{\partial E}{\partial z_{2k}} \cdot g'(a_{2k}) \cdot z_{1j} = S_{2k} \cdot z_{1j}$$

$$\frac{\partial E}{\partial W_{1ij}} = \frac{\partial E}{\partial z_{1j}} \cdot g'(a_{1j}) \cdot x_i = S_{1j} \cdot x_i$$

$$\textcircled{1} \quad S_{2k} = \frac{\partial E}{\partial z_{2k}} \cdot g'(a_{2k})$$

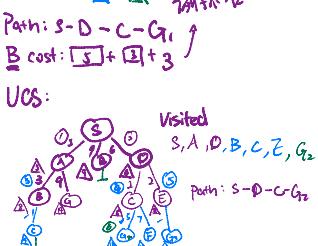
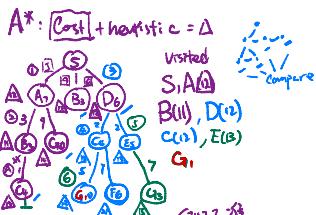
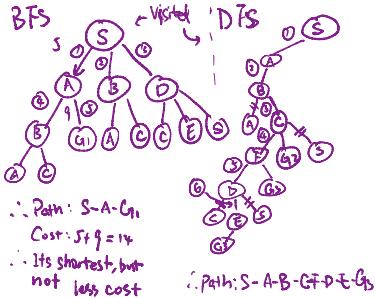
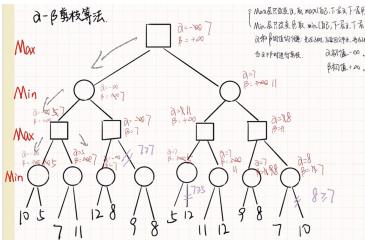
$$\textcircled{2} \quad S_{1j} = \left( \sum_k S_{2k} \cdot W_{2jk} \right) g'(a_{1j}) \cdot x_i = S_{1j} \cdot x_i$$

# Search

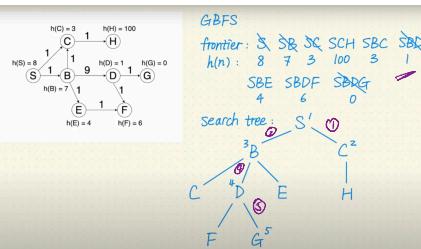


**CONSTRAINT SATISFACTION**

- A set of variables,  $X_1, X_2, X_3 \dots X_n$
- Each variable  $X_i$  has a non-empty domain of possible values, e.g.  $X_1 = \{1, 2, \dots, 4\}$
- A set of constraints, e.g.  $X_1 \neq 3$ ,  $X_3 > X_4$  (using binary)
- Solution: assign a value to each variable such that none of the constraints are broken
- Simplest algorithm: backtracking



Greedy, not optimal or less cost



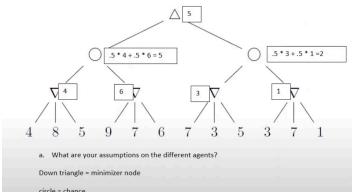
$$\begin{aligned}
 A &= \{1, 2, 3\} & A = 1, B = 1 & \times \\
 B &= \{1, 2, 3\} & A = 1, B = 2 & \times \\
 C &= \{1, 2, 3\} & A = 1, B = 3 & \times \\
 A > B & & A = 2, B = 1, C = 1 & \times \\
 B \neq C & & A = 2, B = 1, C = 2 & \times \\
 A \neq C & & A = 2, B = 1, C = 3 & \checkmark
 \end{aligned}$$

## THE AC-3 ALGORITHM

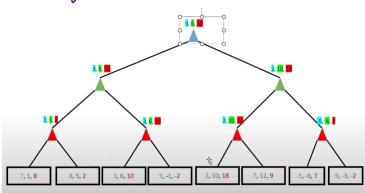
- Turn each binary constraint into two arcs e.g.  $A \neq B$  becomes  $A \neq B$  and  $B \neq A$ .
- Add all the arcs to an agenda
- Repeat until agenda is empty:
  - take an arc  $(X_i, X_j)$  off the agenda
  - for every value of  $X_i$ , there must be some value of  $X_j$
  - remove any inconsistent value from  $X_i$
  - if  $X_i$  has changed, add all arcs of the form  $(X_i, X_k)$  to the agenda

- Evaluation Function  
↳ Produce estimates of Value of non-terminal utilities  
 $Eval(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$   
Eg.  $Eval(s) = 2 \cdot agent\_king(s) + agent\_pawns(s) - 2 \cdot opponent\_king(s) - opponent\_pawns(s)$

## - Expectimax

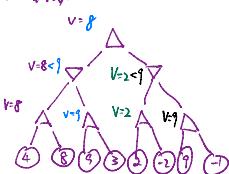


## Multagent



eg.  
 ↓ start  
 Arcs  
 $A \rightarrow B$   
 $B \rightarrow A$  *already in agenda*  
 $B \rightarrow C$   
 $C \rightarrow B$   
 $A \rightarrow B$   
 $B \rightarrow C$   
 $A \rightarrow C$   
 $B \rightarrow A$   
 $C \rightarrow B$   
*nothing change*  
 $B \rightarrow C$

## Minimax



## Minimax - $\alpha$ - $\beta$ Pruning

**MINIMAX WITH  $\alpha$ - $\beta$  PRUNING**

```

MIN-Value ( $s, \alpha, \beta$ ):
  if terminal ( $s$ ) return  $U(s)$ 
   $V = -\infty$ 
  for  $c$  in next-states ( $s$ ):
     $V' = \text{Min-value} (c, \alpha, \beta)$ 
    if  $V' > V$ ,  $V = V'$ 
    if  $V' > \beta$ , return  $V$ 
    if  $V' < \alpha$ ,  $\alpha = V'$ 
  return  $V$ 

MAX-Value ( $s, \alpha, \beta$ ):
  if terminal ( $s$ ) return  $U(s)$ 
   $V = +\infty$ 
  for  $c$  in next-states ( $s$ ):
     $V' = \text{Max-value} (c, \alpha, \beta)$ 
    if  $V' < V$ ,  $V = V'$ 
    if  $V' \leq \alpha$ , return  $V$ 
    if  $V' < \beta$ ,  $\beta = V'$ 
  return  $V$ 

```

