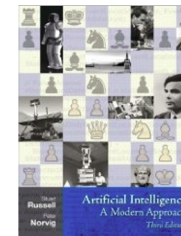# Adversarial Search and Reinforcement Learning
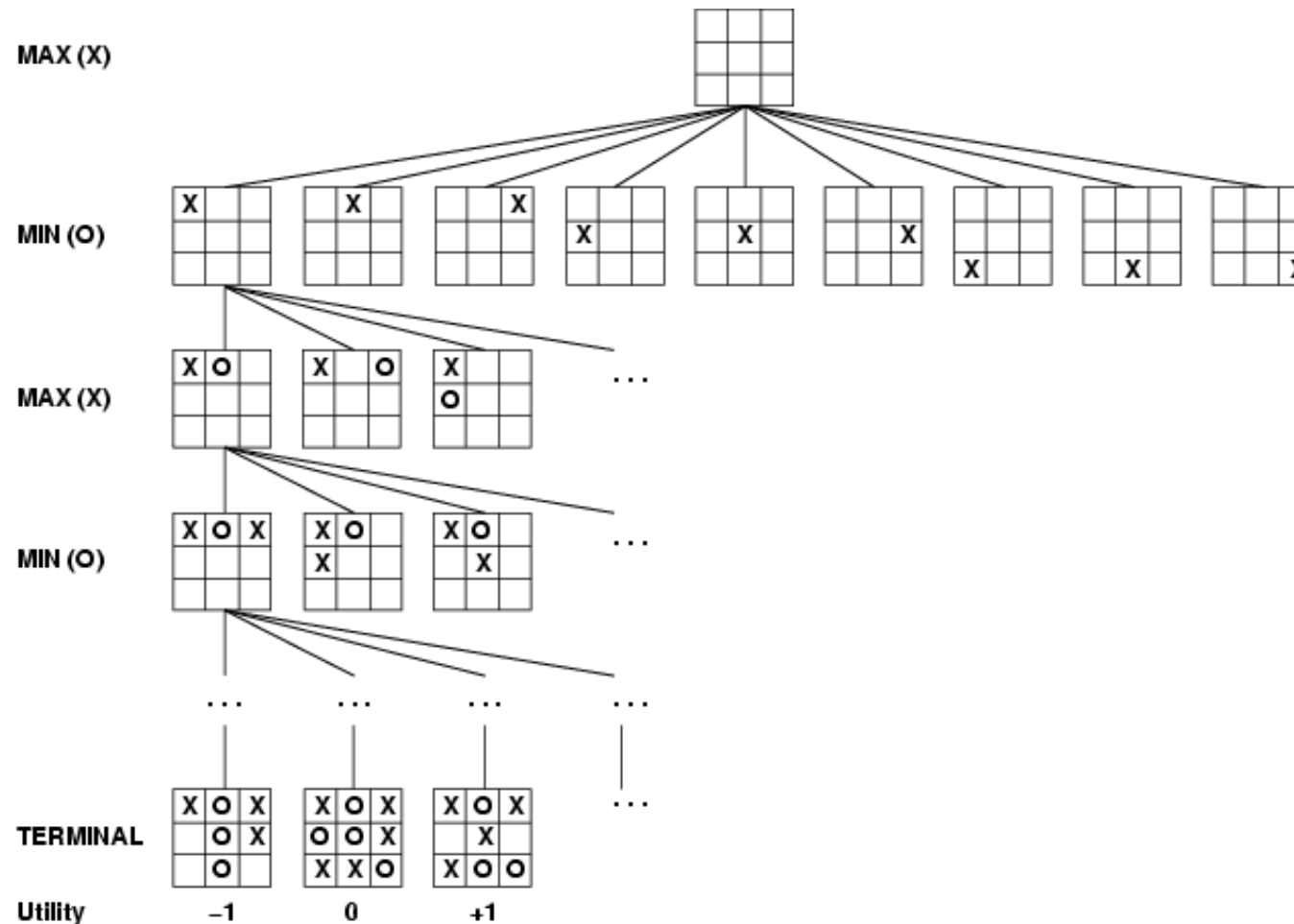
## Chapters 5 & 22

Many slides based on Russell & Norvig's slides
Artificial Intelligence: A Modern Approach

Slides also due to Sriraam Natarajan, Matthew Taylor and Eric Sandholm

# Games vs. search problems

- "Unpredictable" opponent → specifying a move for every possible opponent reply
  - Impossible
  - Unrealistic
- Time limits → unlikely to find goal, must approximate
- Most games are
  - Deterministic
  - Turn-taking
  - Two player
  - Zero-sum games
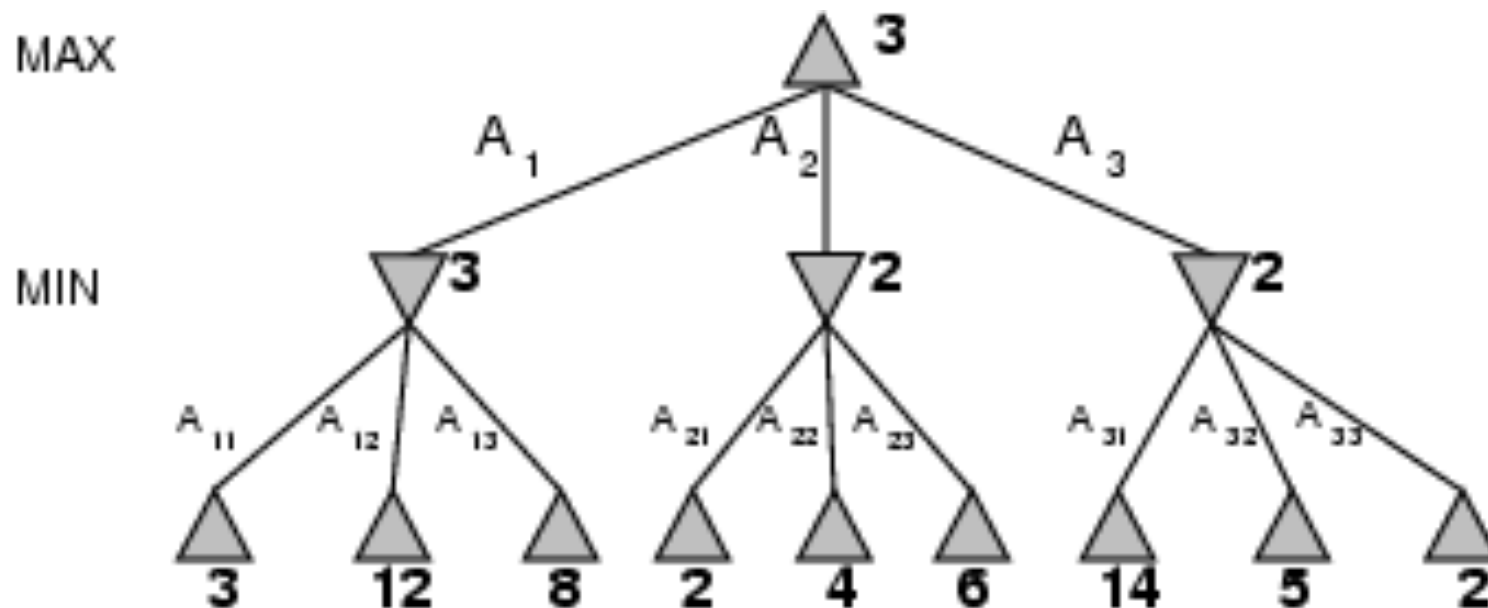- Most real games are stochastic, parallel, multi-agent and utility based

# Game tree (2-player, deterministic, turns)



This is a tic-tac-toe game. The tree is supposedly small. Fewer than 9! = 362,880 terminal nodes. Chess, this number is $10^{40}$.
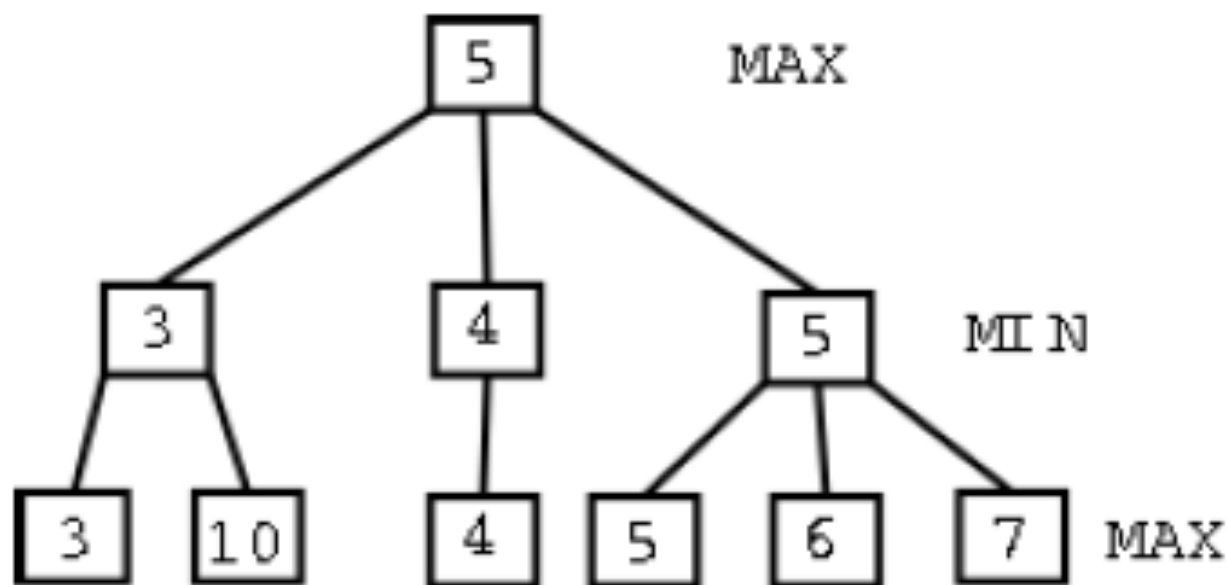
# Minimax

- Perfect play for deterministic games
- Idea: choose move to position with highest minimax value
  = best achievable payoff against best play
- E.g., 2-ply game:

# Example

- You are the "max" player
- Opponent is the min player
- (S)He wants to minimize your score
- You want to maximize
- Key assumption: (S)He is optimal
- How can this be extended to multiple players?
- Replace each utility with a vector of utilities corresponding to each player

# Minimax algorithm

**function** MINIMAX-DECISION(*state*) **returns** *an action*

    $v \leftarrow$ MAX-VALUE(*state*)
    **return** the *action* in SUCCESSORS(*state*) with value $v$

---

**function** MAX-VALUE(*state*) **returns** *a utility value*

    **if** TERMINAL-TEST(*state*) **then return** UTILITY(*state*)
    $v \leftarrow -\infty$
    **for** $a, s$ in SUCCESSORS(*state*) **do**
        $v \leftarrow$ MAX($v$, MIN-VALUE($s$))
    **return** $v$

---

**function** MIN-VALUE(*state*) **returns** *a utility value*

    **if** TERMINAL-TEST(*state*) **then return** UTILITY(*state*)
    $v \leftarrow \infty$
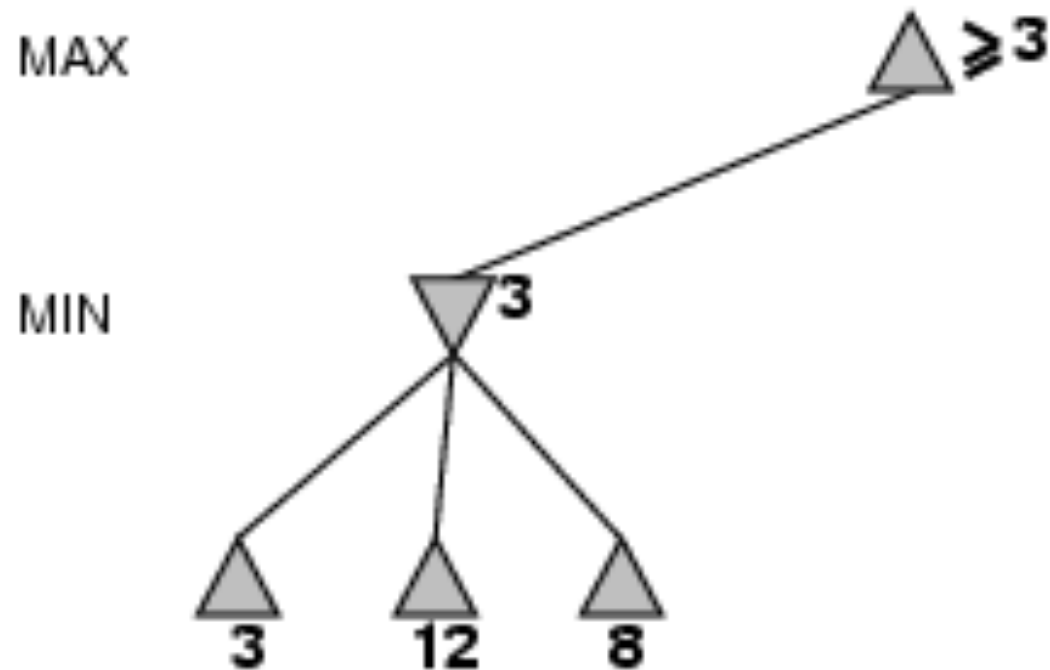    **for** $a, s$ in SUCCESSORS(*state*) **do**
        $v \leftarrow$ MIN($v$, MAX-VALUE($s$))
    **return** $v$
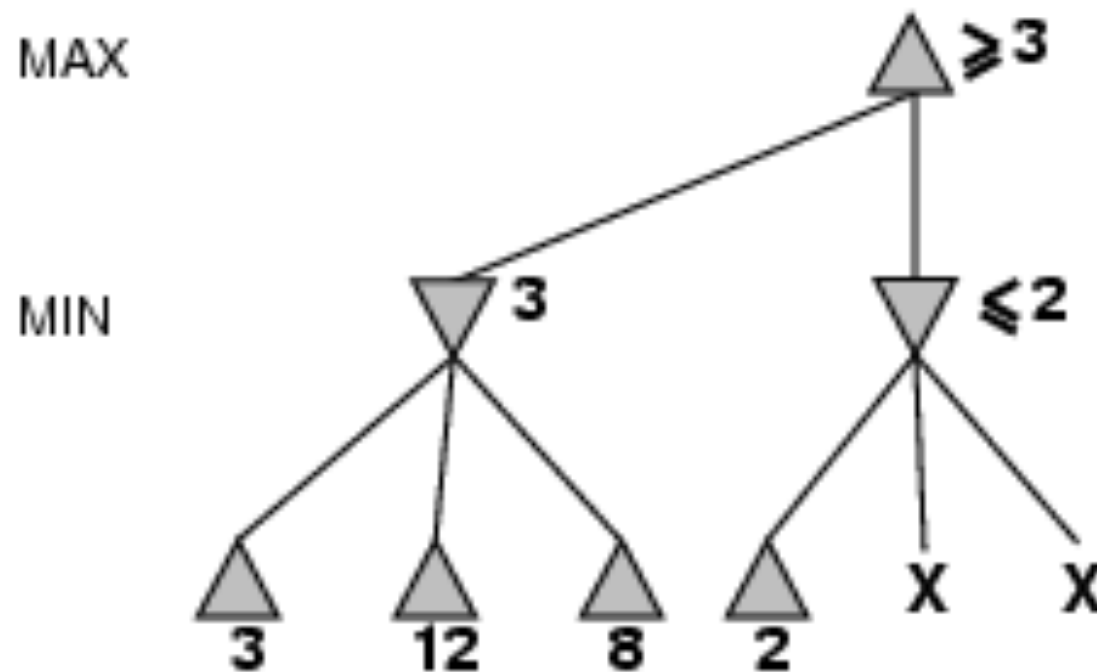
# Properties of minimax

- <u>Complete?</u> Yes (if tree is finite)
- <u>Optimal?</u> Yes (against an optimal opponent)
- <u>Time complexity?</u> $O(b^m)$
- <u>Space complexity?</u> $O(bm)$ (depth-first exploration)
- For chess, $b \approx 35$, $m \approx 100$ for "reasonable" games
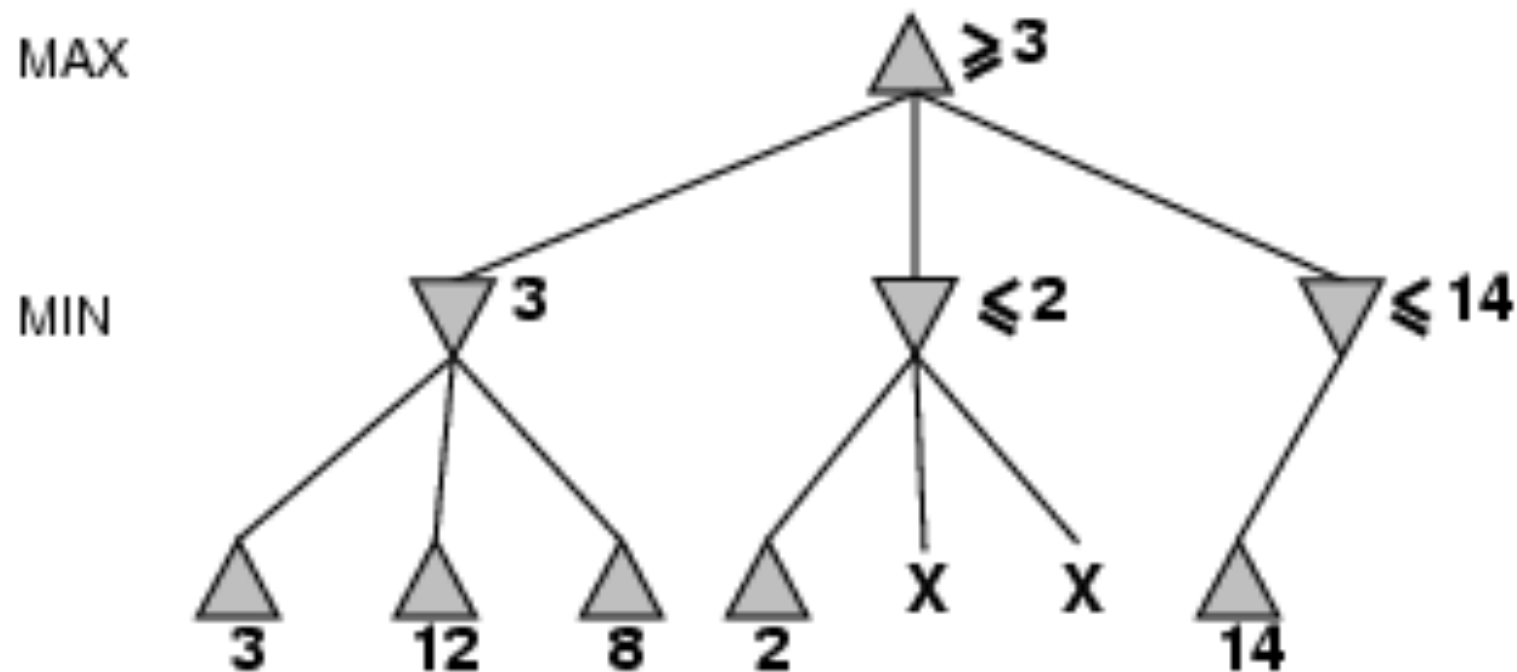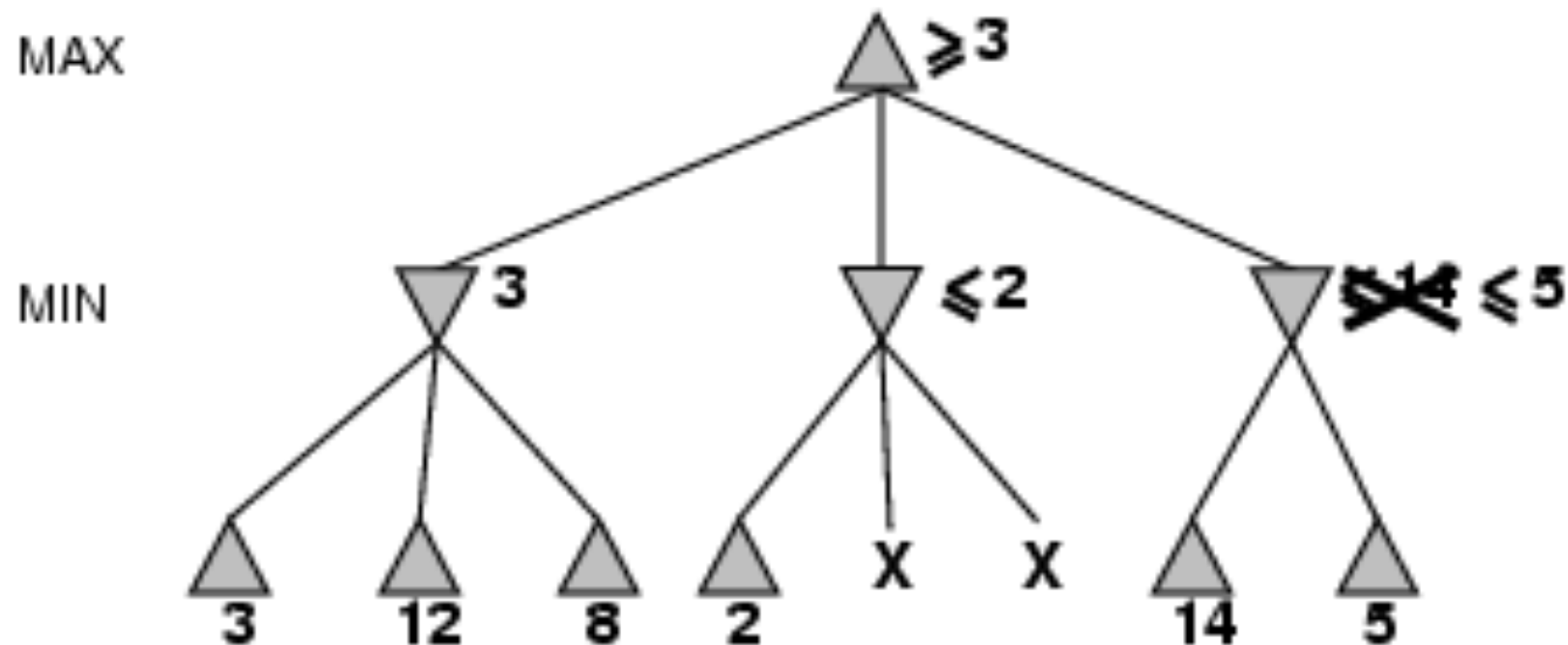  → exact solution completely infeasible
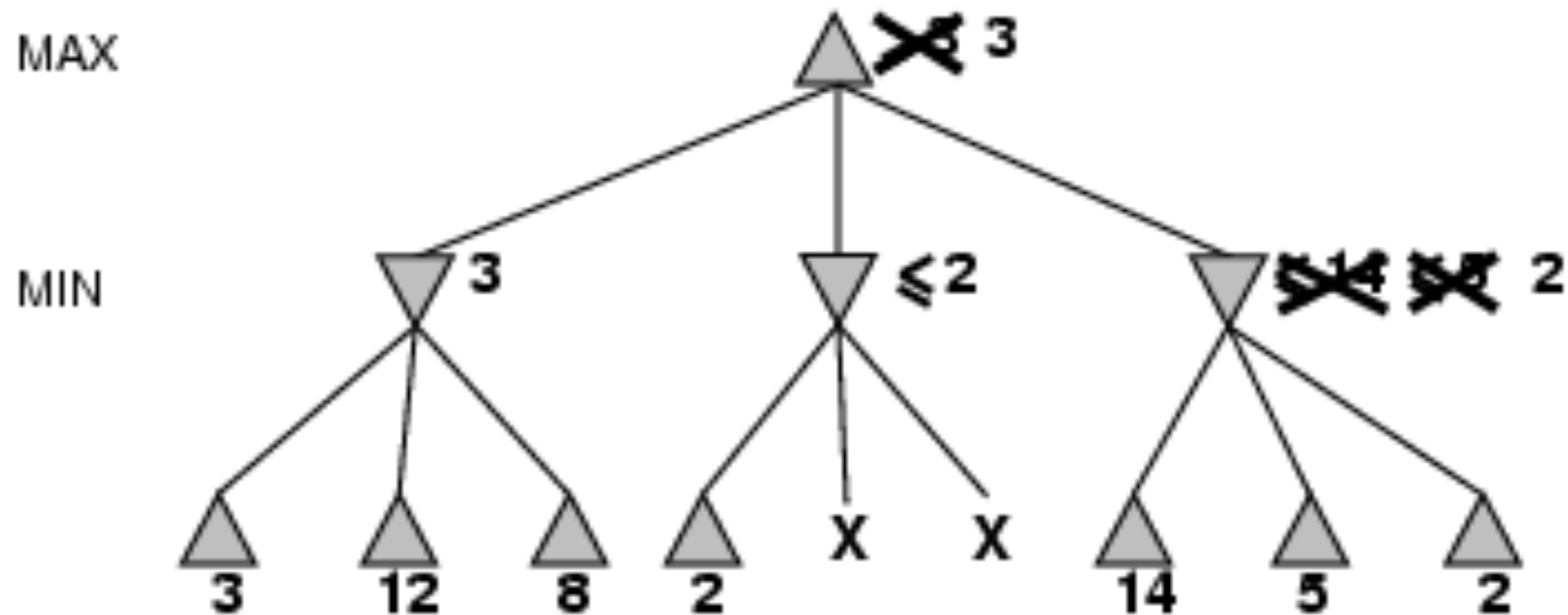
# α-β pruning example

# α-β pruning example
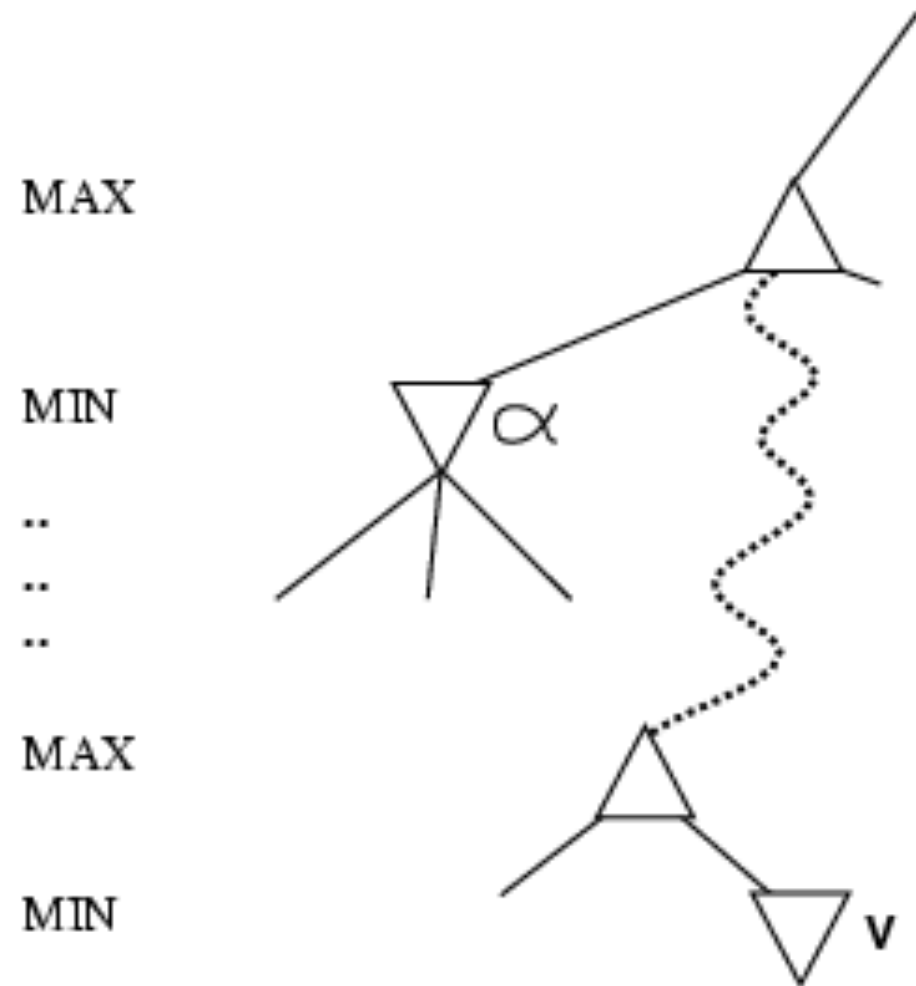
# α-β pruning example

# α-β pruning example

# α-β pruning example

# Why is it called α-β?

- α is the value of the best (i.e., highest-value) choice found so far at any choice point along the path for max

- If v is worse than α, max will avoid it → prune that branch

- Define β similarly for min

# The α-β algorithm

**function** ALPHA-BETA-SEARCH($state$) **returns** $an\ action$
   **inputs**: $state$, current state in game

   $v \leftarrow$ MAX-VALUE($state, -\infty, +\infty$)
   **return** the $action$ in SUCCESSORS($state$) with value $v$

---

**function** MAX-VALUE($state, \alpha, \beta$) **returns** $a\ utility\ value$
   **inputs**: $state$, current state in game
            $\alpha$, the value of the best alternative for MAX along the path to $state$
            $\beta$, the value of the best alternative for MIN along the path to $state$

   **if** TERMINAL-TEST($state$) **then return** UTILITY($state$)
   $v \leftarrow -\infty$
   **for** $a, s$ in SUCCESSORS($state$) **do**
      $v \leftarrow$ MAX($v$, MIN-VALUE($s, \alpha, \beta$))
      **if** $v \geq \beta$ **then return** $v$
      $\alpha \leftarrow$ MAX($\alpha, v$)
   **return** $v$

# The α-β algorithm

**function** MIN-VALUE($state, \alpha, \beta$) **returns** *a utility value*
    **inputs**: *state*, current state in game
                  $\alpha$, the value of the best alternative for MAX along the path to *state*
                  $\beta$, the value of the best alternative for MIN along the path to *state*

    **if** TERMINAL-TEST($state$) **then return** UTILITY($state$)
    $v \leftarrow +\infty$
    **for** $a, s$ **in** SUCCESSORS($state$) **do**
        $v \leftarrow$ MIN($v$, MAX-VALUE($s, \alpha, \beta$))
        **if** $v \leq \alpha$ **then return** $v$
        $\beta \leftarrow$ MIN($\beta, v$)
    **return** $v$

# Evaluation functions

- For chess, typically <span style="color:red">linear</span> weighted sum of <span style="color:blue">features</span>

$$Eval(s) = w_1 \, f_1(s) + w_2 \, f_2(s) + \ldots + w_n \, f_n(s)$$

- e.g., $w_1 = 9$ with

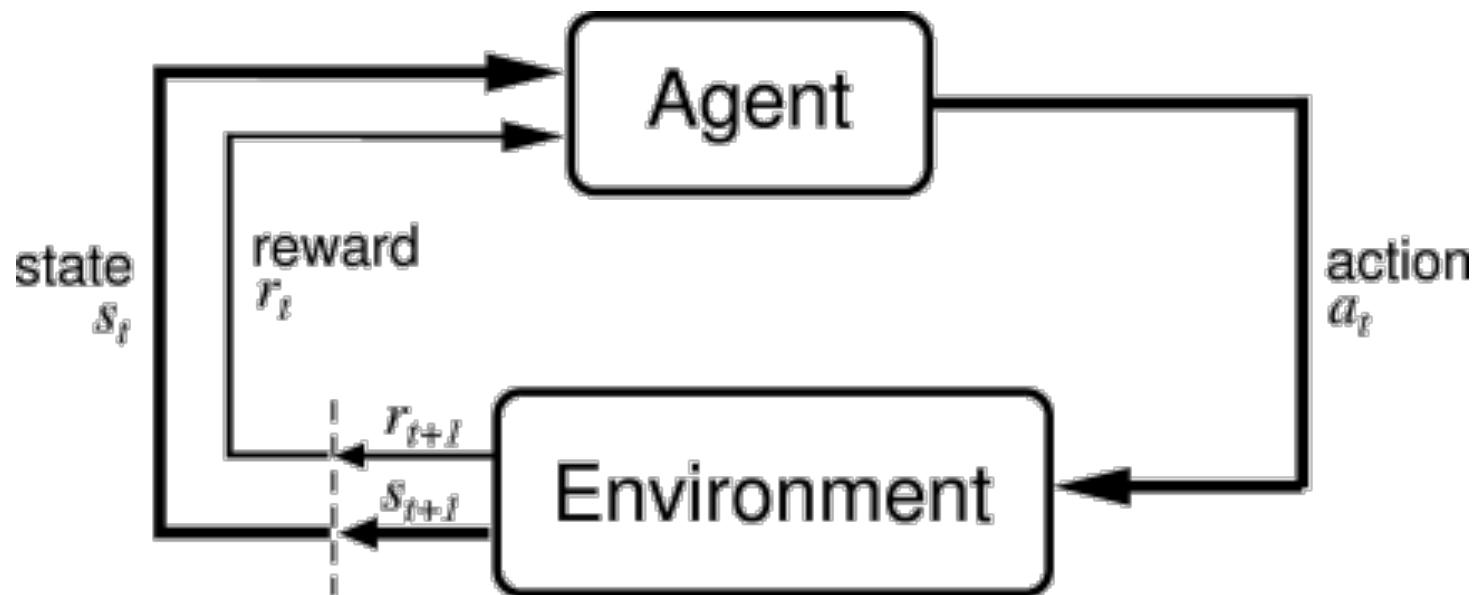  $f_1(s)$ = (number of white queens) – (number of black queens), etc.

# Deterministic games in practice

- Checkers: Chinook ended 40-year-reign of human world champion Marion Tinsley in 1994. Used a precomputed endgame database defining perfect play for all positions involving 8 or fewer pieces on the board, a total of 444 billion positions.

- Chess: Deep Blue defeated human world champion Garry Kasparov in a six-game match in 1997. Deep Blue searches 200 million positions per second, uses very sophisticated evaluation, and undisclosed methods for extending some lines of search up to 40 ply.

- Othello: human champions refuse to compete against computers, who are too good.

- Go: human champions refused to compete against computers, who were supposedly too bad. In go, $b > 300$, so most programs use pattern knowledge bases to suggest plausible moves.

Speaking of Go, we saw already that Deep Networks within Reinforcement Learning may help. But what is Reinforcement Learning actually?

# Reinforcement Learning

- ## Basic idea:
  - Receive feedback in the form of <span style="color:red">rewards</span>
  - Agent's utility is defined by the reward function
  - Must (learn to) act so as to <span style="color:red">maximize expected rewards</span>
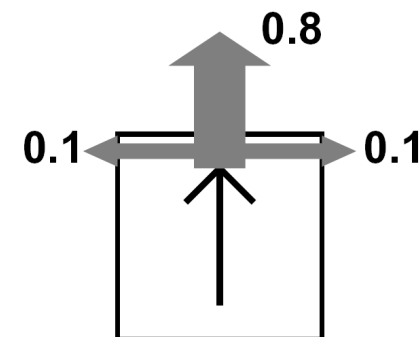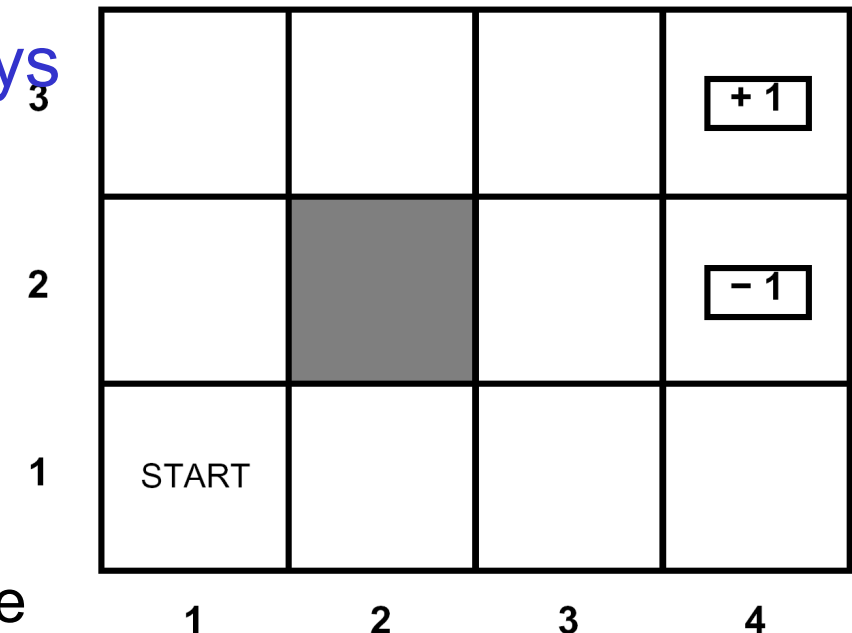
# Reinforcement Learning (RL)

- RL algorithms attempt to find a policy for maximizing cumulative reward for the agent over the course of the problem.

- Typically represented by a **Markov Decision Process**

- RL differs from supervised learning in that correct input/output pairs are never presented, nor sub-optimal actions explicitly corrected.
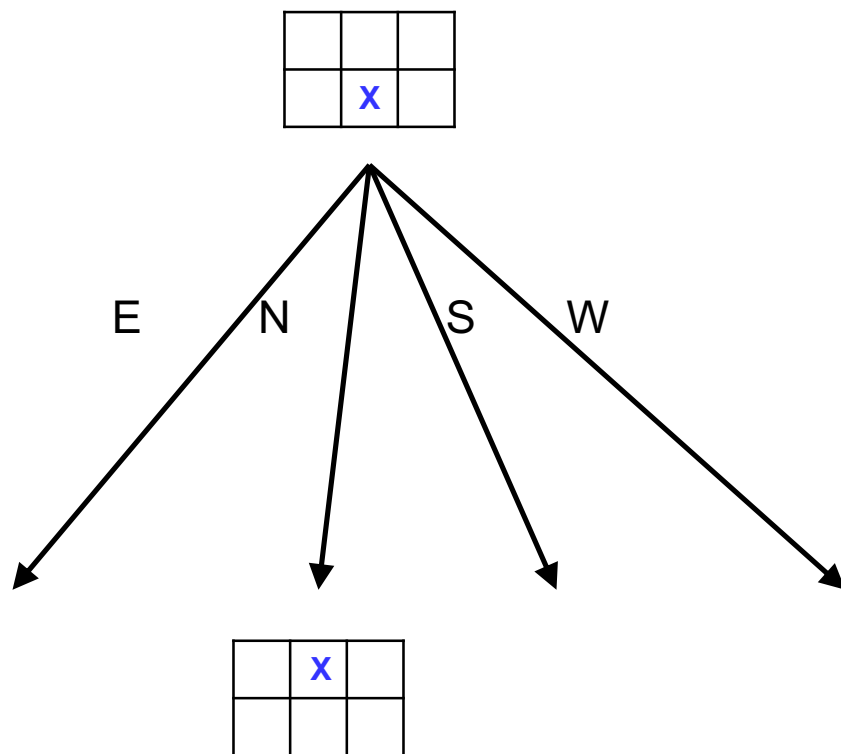
# Grid World

- ## The agent lives in a grid

- ## Walls block the agent's path

- ## The agent's actions do not always go as planned:

  - 80% of the time, the action North takes the agent North
  (if there is no wall there)

  - 10% of the time, North takes the agent West; 10% East

  - If there is a wall in the direction the agent would have been taken, the agent stays put

- ## Small "living" reward each step

- ## Big rewards come at the end

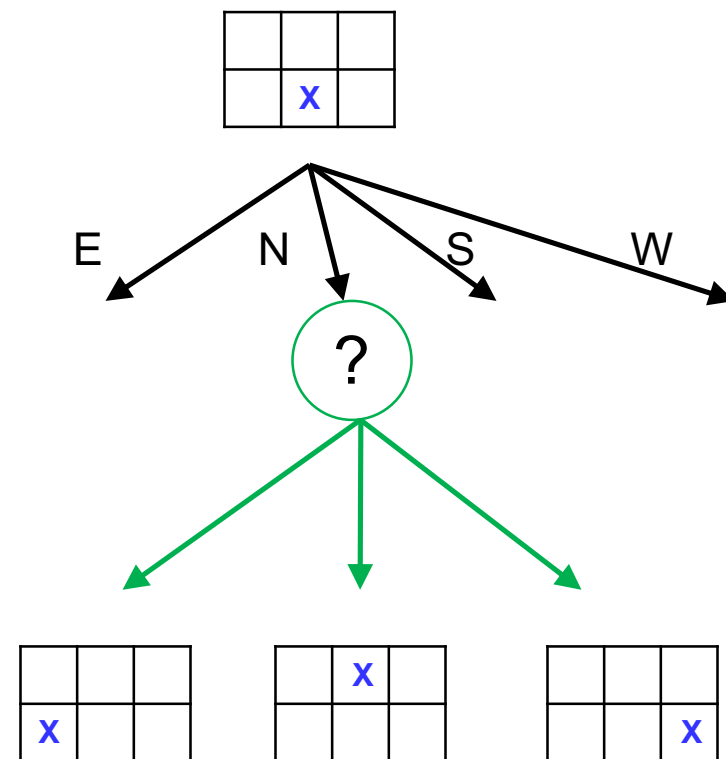- ## Goal: maximize sum of rewards*
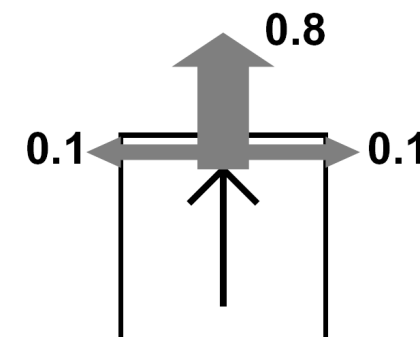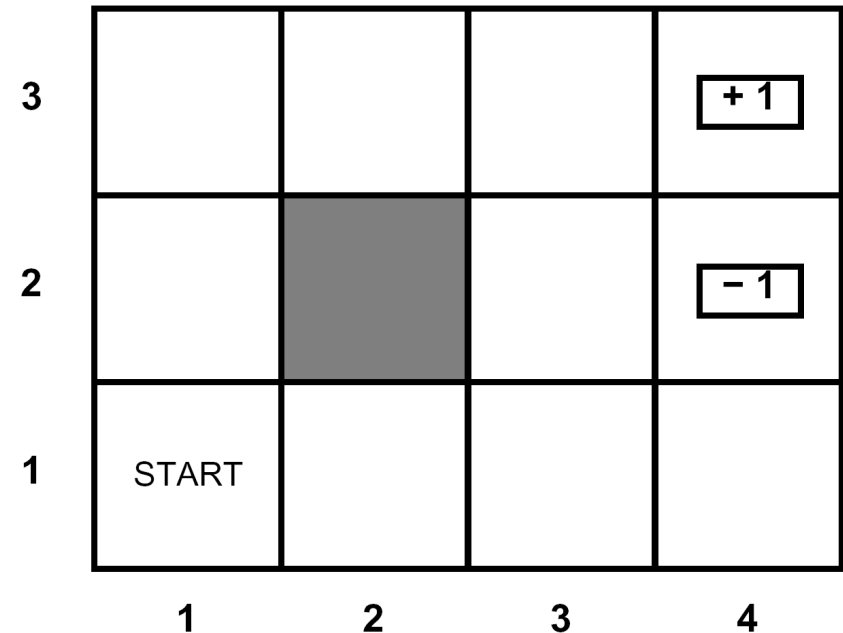
# Grid Futures

## Deterministic Grid World



## Stochastic Grid World

# Markov Decision Processes (MDP)

- An MDP is defined by:
  - A set of states s ∈ S
  - A set of actions a ∈ A
  - A transition function T(s,a,s')
    - Prob that a from s leads to s'
    - i.e., P(s' | s,a)
    - Also called the model
  - A reward function R(s, a, s')
    - Sometimes just R(s) or R(s')
  - A start state (or distribution)
  - Maybe a terminal state
  - A discount factor: γ

- MDPs are a family of non-deterministic search problems

  - Reinforcement learning: MDPs where we don't know the transition or reward functions

# What is Markov about MDPs?

- Andrey Markov (1856-1922)

- "Markov" generally means that given the present state, the future and the past are independent

- For Markov decision processes, "Markov" means:

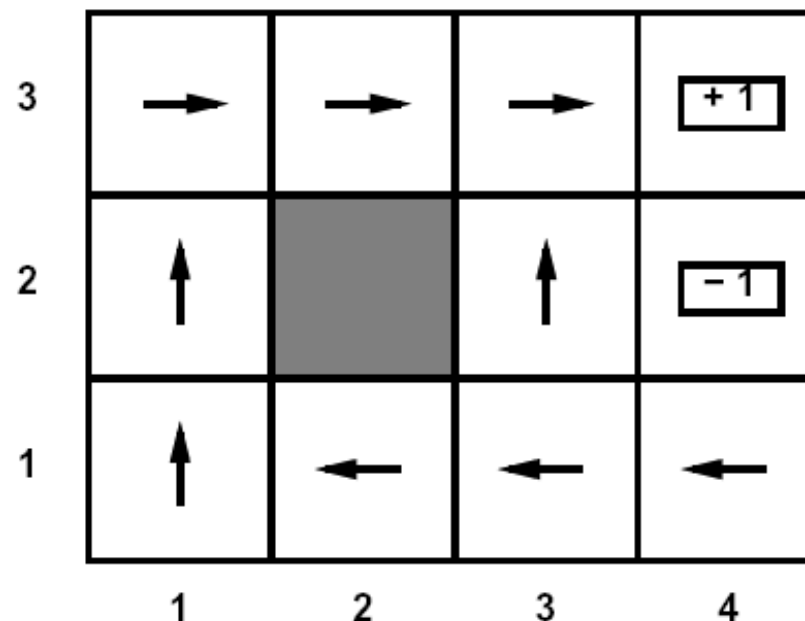$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \ldots S_0 = s_0)$$

$$=$$

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

# Solving MDPs

- ## In deterministic single-agent search problems, want an optimal plan, or sequence of actions, from start to a goal

- ## In an MDP, we want an optimal policy $\pi^*: S \rightarrow A$

  - ### A policy $\pi$ gives an action for each state
  - ### An optimal policy maximizes expected utility if followed
  - ### Defines a reflex agent

Optimal policy when R(s, a, s') = -0.03 for all non-terminals s

# Example Optimal Policies



R(s) = -0.01

R(s) = -0.03

R(s) = -0.4

R(s) = -2.0

# Recap: Defining MDPs

- ## Markov decision processes:
  - ### States S
  - ### Start state $s_0$
  - ### Actions A
  - ### Transitions P(s'|s,a) (or T(s,a,s'))
  - ### Rewards R(s,a,s') (and discount $\gamma$)

- ## MDP quantities so far:
  - ### Policy = Choice of action for each state
  - ### Utility (or return) = sum of discounted rewards

# Optimal Utilities

- Fundamental operation: compute the values (optimal expectimax utilities) of states s

- Why?  Optimal values define optimal policies!

- Define the value of a state s:
  $V^*(s)$ = expected utility starting in s and acting optimally
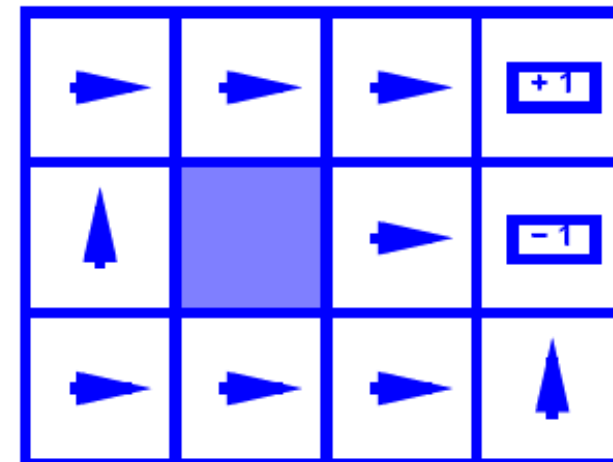
- Define the value of a q-state (s,a):
  $Q^*(s,a)$ = expected utility starting in s, taking action a and thereafter acting optimally

- Define the optimal policy:
  $\pi^*(s)$ = optimal action from state s

# Value Iteration

- ## Idea:
  - Start with $V_0^*(s) = 0$
  - Given $V_i^*$, calculate the values for all states for depth i+1:

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_i(s') \right]$$

  - This is called a value update or Bellman update
  - Repeat until convergence

- ## Theorem: will converge to unique optimal values
  - Basic idea: approximations get refined towards optimal values
  - Policy may converge long before values do

# Example: Value Iteration

$V_1$                                                                     $V_2$



- Information propagates outward from terminal states and eventually all states have correct value estimates

# Example: Value Iteration



- Information propagates outward from terminal states and eventually all states have correct value estimates

# Discounted Rewards

- **Rewards in the future are worth less than an immediate reward**

- Discount factor $\gamma \leq 1$     (often $\gamma = 0.9$)
- Assume reward $n$ years in the future is only worth $(\gamma)^n$ of the value of immediate reward

  - $(0.9\text{^}6) * 10{,}000 = 0.531 * 10{,}000 = 5310$

- For each state, calculate a *utility* value equal to the *Sum of Future Discounted Rewards*

# Reinforcement Learning

- Reinforcement learning:
  - Still assume an MDP:
    - A set of states s ∈ S
    - A set of actions (per state) A
    - A model T(s,a,s')
    - A reward function R(s,a,s')
    - A discount factor $\gamma$ (could be 1)
  - Still looking for a policy $\pi$(s)

  - New twist: don't know T or R
    - i.e. don't know which states are good or what the actions do
    - Must actually try actions and states out to learn

# Passive Learning

- ## Simplified task
  - ### You don't know the transitions $T(s,a,s')$
  - ### You don't know the rewards $R(s,a,s')$
  - ### You are given a policy $\pi(s)$
  - ### Goal: learn the state values
  - ### … what policy evaluation did

- ## In this case:
  - ### Learner "along for the ride"
  - ### No choice about what actions to take
  - ### Just execute the policy and learn from experience
  - ### We'll get to the active case soon
  - ### This is NOT offline planning!  You actually take actions in the world and see what happens…

# Model-Based Learning

- Idea:
  - Learn the model empirically through experience
  - Solve for values as if the learned model were correct

- Simple empirical model learning
  - Count outcomes for each s,a
  - Normalize to give estimate of **T(s,a,s')**
  - Discover **R(s,a,s')** when we experience (s,a,s')

- Solving the MDP with the learned model
  - Iterative policy evaluation, for example

$$V_{i+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_i^{\pi}(s')]$$

# Sample-Based Policy Evaluation?

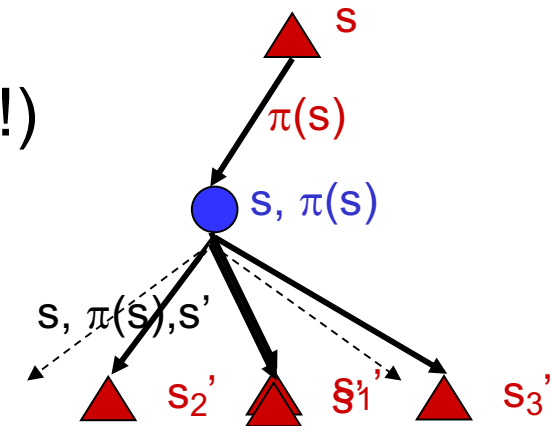$$V_{i+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_i^{\pi}(s')]$$

- Who needs T and R?  Approximate the expectation with samples (drawn from T!)

$$sample_1 = R(s, \pi(s), s_1') + \gamma V_i^{\pi}(s_1')$$

$$sample_2 = R(s, \pi(s), s_2') + \gamma V_i^{\pi}(s_2')$$

$$\ldots$$

$$sample_k = R(s, \pi(s), s_k') + \gamma V_i^{\pi}(s_k')$$

s

π(s)

s, π(s)

s, π(s),s'
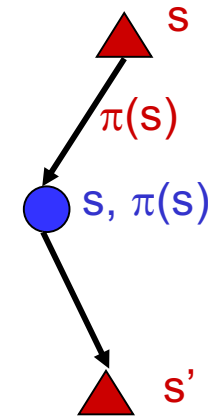
$s_2$'    $s_1$'    $s_3$'

$$V_{i+1}^{\pi}(s) \leftarrow \frac{1}{k} \sum_i sample_i$$

*Almost!  But we only actually make progress when we move to i+1.*

# Temporal-Difference Learning

- Big idea: learn from every experience!
  - Update V(s) each time we experience (s,a,s',r)
  - Likely s' will contribute updates more often

- Temporal difference learning
  - Policy still fixed!
  - Move values toward value of whatever successor occurs: running average!

s

$\pi(s)$

s, $\pi(s)$

s'

**Sample of V(s):**   $$sample = R(s, \pi(s), s') + \gamma V^\pi(s')$$

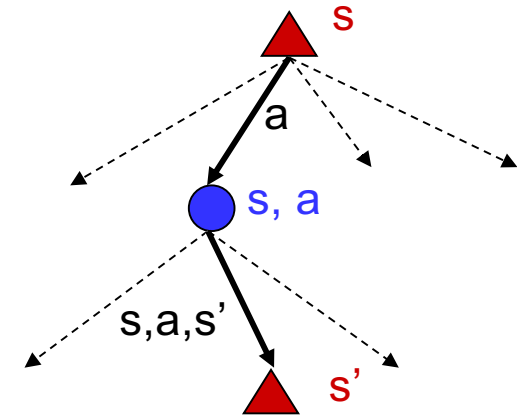**Update to V(s):**   $$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$$

**Same update:**   $$V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$$

# Problems with TD Value Learning

- TD value leaning is a model-free way to do policy evaluation

- However, if we want to turn values into a (new) policy, we're sunk:
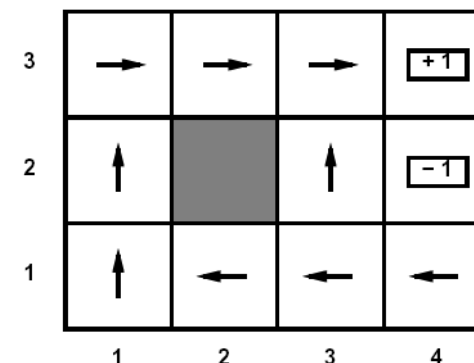
$$\pi(s) = \arg\max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

- Idea: learn Q-values directly

- Makes action selection model-free too!

# Active Learning

- ## Full reinforcement learning
  - You don't know the transitions T(s,a,s')
  - You don't know the rewards R(s,a,s')
  - You can choose any actions you like
  - Goal: learn the optimal policy
  - … what value iteration did!

- ## In this case:
  - Learner makes choices!
  - Fundamental tradeoff: exploration vs. exploitation
  - This is NOT offline planning! You actually take actions in the world and find out what happens…

# Q-Learning

- Q-Learning: sample-based Q-value iteration
- Learn Q*(s,a) values
  - Receive a sample (s,a,s',r)
  - Consider your old estimate: $Q(s,a)$
  - Consider your new sample estimate:

$$Q^*(s,a) = \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q^*(s',a') \right]$$

$$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$

  - Incorporate the new estimate into a running average:

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha)[sample]$$

# Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy
  - If you explore enough
  - If you make the learning rate small enough
  - … but not decrease it too quickly!
  - Basically doesn't matter how you select actions (!)

# **Exploration / Exploitation**

- Several schemes for forcing exploration
  - Simplest: random actions ($\varepsilon$ greedy)
    - Every time step, flip a coin
    - With probability $\varepsilon$, act randomly
    - With probability 1-$\varepsilon$, act according to current policy

  - Problems with random actions?
    - You do explore the space, but keep thrashing around once learning is done
    - One solution: lower $\varepsilon$ over time
    - Another solution: exploration functions

# The Story So Far: MDPs and RL

## Things we know how to do:

- ## If we know the MDP
  - Compute V*, Q*, $\pi$* exactly
  - Evaluate a fixed policy $\pi$

- ## If we don't know the MDP
  - We can estimate the MDP then solve

  - We can estimate V for a fixed policy $\pi$
  - We can estimate Q*(s,a) for the optimal policy while executing an exploration policy
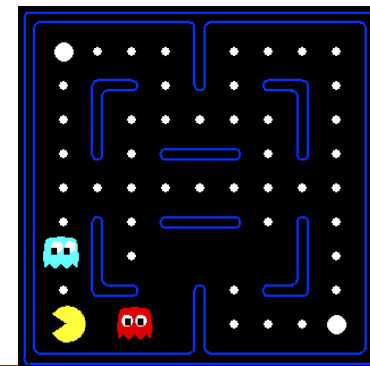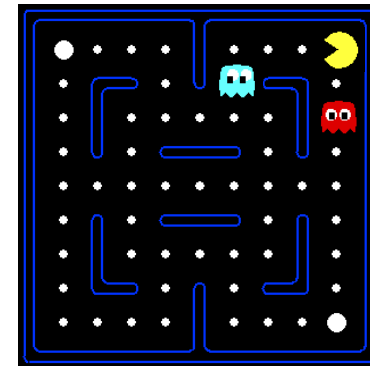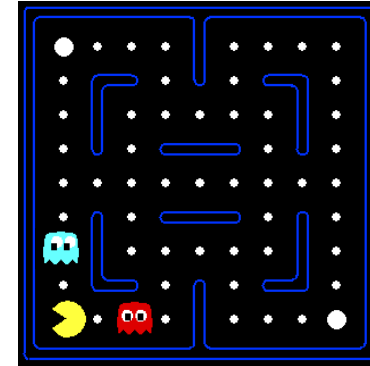
## Techniques:

- ## Model-based DPs
  - Value and policy Iteration
  - Policy evaluation

- ## Model-based RL

- ## Model-free RL:
  - Value learning
  - Q-learning

# Q-Learning

- In realistic situations, we cannot possibly learn about every single state!
  - Too many states to visit them all in training
  - Too many states to hold the q-tables in memory

- Instead, we want to generalize:
  - Learn about some small number of training states from experience
  - Generalize that experience to new, similar states
  - This is a fundamental idea in machine learning, and we'll see it over and over again
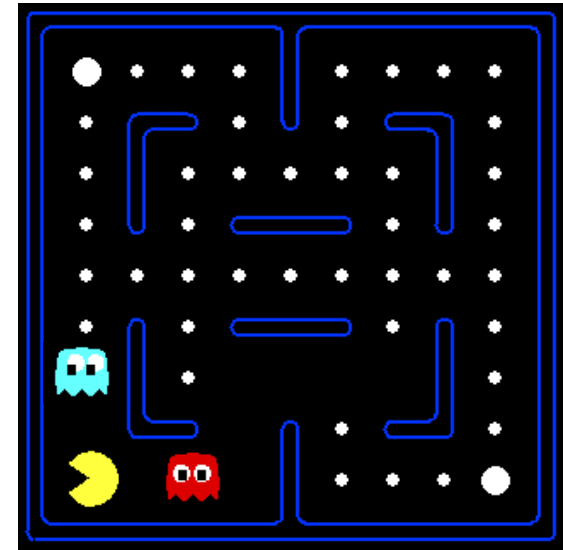
# Example: Pacman

- Let's say we discover through experience that this state is bad:



- In naïve q learning, we know nothing about this state or its q states:



- Or even this one!

# Feature-Based Representations

- Solution: describe a state using a vector of features
  - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - Distance to closest dot
    - Number of ghosts
    - $1 / (\text{dist to dot})^2$
    - Is Pacman in a tunnel? (0/1)
    - …… etc.
  - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)

# Linear Feature Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \ldots + w_n f_n(s)$$

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but be very different in value!

# Function Approximation

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

- Q-learning with linear q-functions:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \, [error]$$

$$w_i \leftarrow w_i + \alpha \, [error] \, f_i(s,a)$$

- Intuitive interpretation:
  - Adjust weights of active features
  - E.g. if something unexpectedly bad happens, disprefer all states with that state's features

- Formal justification: online least squares

# Example: Q-Pacman

$$Q(s,a) = 4.0 f_{DOT}(s,a) - 1.0 f_{GST}(s,a)$$

$$f_{DOT}(s, \text{NORTH}) = 0.5$$
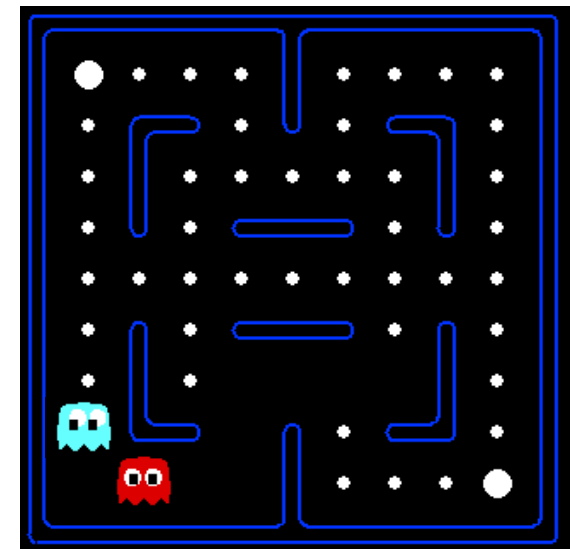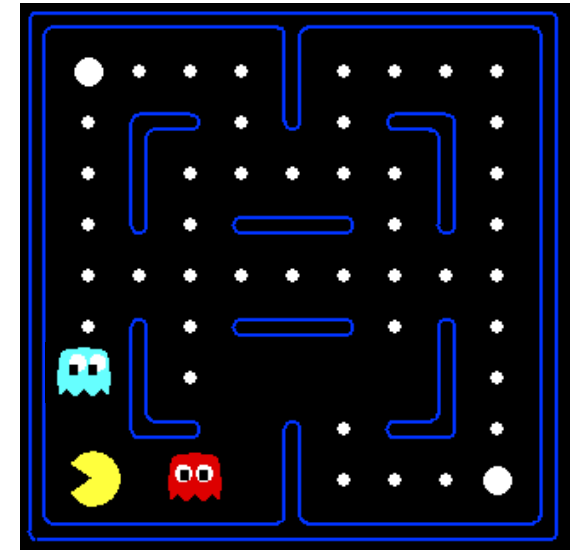
$$f_{GST}(s, \text{NORTH}) = 1.0$$

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

$$R(s,a,s') = -500$$

$$error = -501$$

$$w_{DOT} \leftarrow 4.0 + \alpha\,[-501]\,0.5$$

$$w_{GST} \leftarrow -1.0 + \alpha\,[-501]\,1.0$$

$$Q(s,a) = 3.0 f_{DOT}(s,a) - 3.0 f_{GST}(s,a)$$

# Policy Search



http://heli.stanford.edu/

# Policy Search

- Problem: often the feature-based policies that work well aren't the ones that approximate V / Q best

- Solution: learn the policy that maximizes rewards rather than the value that predicts rewards

- This is the idea behind policy search, such as what controlled the upside-down helicopter.

- Genetic Algorithms can be used as a type of policy search.

# Policy Search

- ## Simplest policy search:
  - Start with an initial linear value function or q-function
  - Nudge each feature weight up and down and see if your policy is better than before

- ## Problems:
  - How do we tell the policy got better?
  - Need to run many sample episodes!
  - If there are a lot of features, this can be impractical