# Recap

.

# Recap

DMU

# Recap

DMU

Probabilistic Models

# Recap

DMU

- Probabilistic Models
- MDPs
- Reinforcement Learning
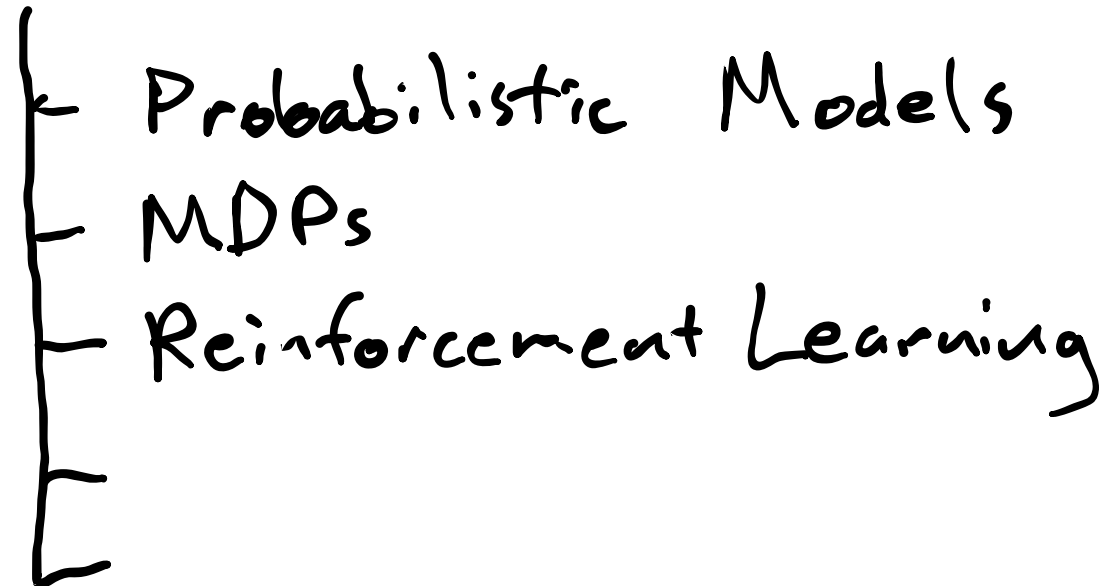
# Recap

DMU
- Probabilistic Models
- MDPs
- Reinforcement Learning
- POMDPs
- Games

# Probabilistic Models

# Probabilistic Models

**3 Rules**

$P(A)$

$P(A, B)$

$P(A | B)$

# Probabilistic Models

**3 Rules**

$P(A)$

$P(A,B)$

$P(A \mid B)$

1. $0 \leq P(X \mid Y) \leq 1$

$\sum_{x \in X} P(x \mid Y) = 1$

# Probabilistic Models

**3 Rules**

$P(A)$

$P(A, B)$

$P(A \mid B)$

1. $0 \leq P(X \mid Y) \leq 1$
   $\sum_{x \in X} P(x \mid Y) = 1$

2. $P(X) = \sum_{y \in Y} P(X, y)$

# Probabilistic Models

$P(A)$

$P(A,B)$

$P(A|B)$

**3 Rules**

1. $0 \leq P(X \mid Y) \leq 1$
   $\sum_{x \in X} P(x \mid Y) = 1$

2. $P(X) = \sum_{y \in Y} P(X, y)$

3. $P(X \mid Y) = \frac{P(X,Y)}{P(Y)}$

# Probabilistic Models

## 3 Rules

$P(A)$

$P(A, B)$

$P(A | B)$

1. $0 \leq P(X \mid Y) \leq 1$
   $\sum_{x \in X} P(x \mid Y) = 1$

2. $P(X) = \sum_{y \in Y} P(X, y)$

3. $P(X \mid Y) = \frac{P(X,Y)}{P(Y)}$

## Bayes Rule

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

# Probabilistic Models

$P(A)$

$P(A, B)$

$P(A \mid B)$

## 3 Rules

1. $0 \leq P(X \mid Y) \leq 1$
   $\sum_{x \in X} P(x \mid Y) = 1$

2. $P(X) = \sum_{y \in Y} P(X, y)$

3. $P(X \mid Y) = \frac{P(X,Y)}{P(Y)}$

## Bayes Rule
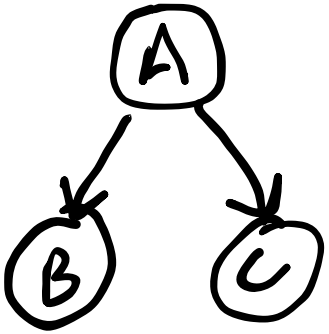
$P(A \mid B) = \dfrac{P(B \mid A)P(A)}{P(B)}$
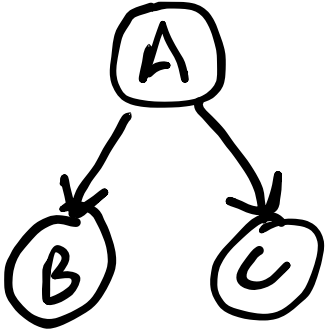
## Independence

$A \perp B \iff P(A, B) = P(A)P(B)$

$A \perp B \mid C \iff P(A, B \mid C) = P(A \mid C)P(B \mid C)$
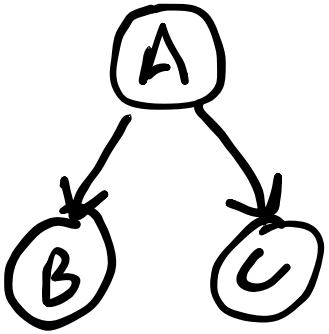
# Bayesian Networks

# Bayesian Networks

# Bayesian Networks



**Chain Rule**

$$P(X_{1:n}) = \prod_i P(X_i \mid Pa(X_i))$$
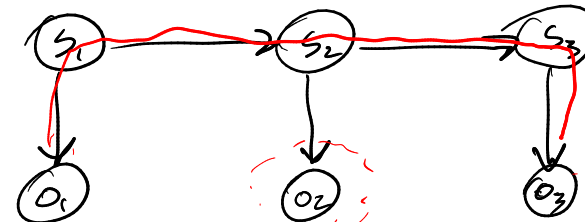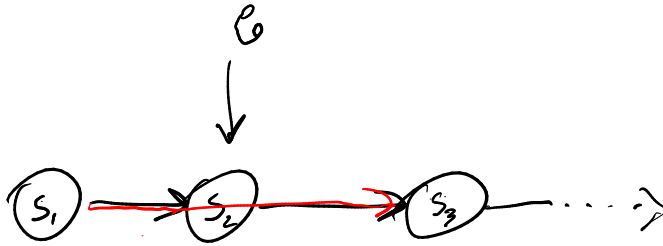
# Bayesian Networks



## Chain Rule

$$P(X_{1:n}) = \prod_i P(X_i \mid Pa(X_i))$$

## Conditional Independence

$X \perp Y \mid \mathcal{C}$ if all paths between $X$ and $Y$ are d-separated by $\mathcal{C}$



$O_1 \perp O_3 \mid O_2$ ?

We cannot prove based on structure

# Markov Decision Processes
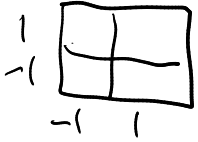
# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Cartesian Product of two sets

$$\{-1, 1\} \times \{-1, 1\} = \{(-1,-1), (-1,1), (1,-1), (1,1)\}$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$s = (x^1, y^1, x^2, y^2, b)$

$b \in \{1, 2\}$

$$S = [-1, 1]^4 \times \{1, 2\}$$

A state is usually represented as a vector or tuple of state variables

A state space is a set of all possible states

$s \in S$

$$s = (x^1, y^1, x^2, y^2)$$

Case 1
$x \in \{-1, 1\}$
$y \in \{-1, 1\}$

$$S = \{-1, 1\} \times \{-1, 1\} \times \{-1, 1\} \times \{-1, 1\} = \{-1, 1\}^4$$

Case 2
$x \in [-1, 1]$
$y \in [-1, 1]$

$$S = [-1, 1] \times [-1, 1] \times [-1, 1] \times [-1, 1] = [-1, 1]^4$$

Case 3
$x \in \{-10, -9, \dots, 9, 10\}$

$$S = \{-10, \dots, 10\}^4$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$$\text{maximize}_\pi \; E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$$\underset{\pi}{\text{maximize}} \quad E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)\right]$$

$$Q^{\pi}(s, a) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a, a_t = \pi(s_t)\right]$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$$\text{maximize}_{\pi} \; E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)\right]$$

$$Q^{\pi}(s, a) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \,\Big|\, s_0 = s, a_0 = a, a_t = \pi(s_t)\right]$$

$$V^{\pi}(s) = Q^{\pi}(s, \pi(s))$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$$\text{maximize}_{\pi} \ E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)\right]$$

$$Q^{\pi}(s, a) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a, a_t = \pi(s_t)\right]$$

$$V^{\pi}(s) = Q^{\pi}(s, \pi(s))$$

$$V^{\pi}(s) = R(s, a) + \gamma E[V^{\pi}(s')]$$

$$V^*(s) = \max_a \{R(s, a) + \gamma E[V^*(s')]\}$$

$$B[V](s) = \max_a \{R(s, a) + \gamma E[V(s')]\}$$

# Markov Decision Processes

$$(S, A, R, T, \gamma)$$

Examples: $S = \{1, 2, 3\}$ or $S = \mathbb{R}^2$

$$s = (x, \dot{x}) \in S = \mathbb{R}^2$$

$$\underset{\pi}{\text{maximize}} \ E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

$$Q^\pi(s, a) = E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a, a_t = \pi(s_t) \right]$$

$$V^\pi(s) = Q^\pi(s, \pi(s))$$

$$V^\pi(s) = R(s, a) + \gamma E[V^\pi(s')] \qquad \text{Policy Evaluation}$$

$$V^*(s) = \max_a \{R(s, a) + \gamma E[V^*(s')]\} \qquad \text{Bellman's Equation: Certificate of Optimality}$$

$$B[V](s) = \max_a \{R(s, a) + \gamma E[V(s')]\} \qquad \text{Bellman's Operator}$$

# Offline MDP Algorithms

# Offline MDP Algorithms

**Policy Iteration**

loop
    Evaluate Policy
    Improve Policy

# Offline MDP Algorithms

**Policy Iteration**

loop
    Evaluate Policy
    Improve Policy

**Value Iteration**

loop
    $V \leftarrow B[V]$

# Offline MDP Algorithms

**Policy Iteration**

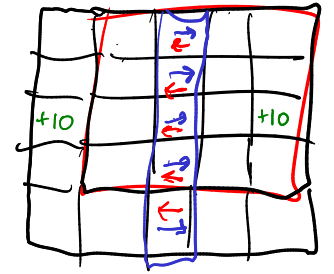loop

    Evaluate Policy

    Improve Policy

Converges because
policy always improves
and there are a finite
number of policies

**Value Iteration**

loop

    $V \leftarrow B[V]$

# Offline MDP Algorithms

**Policy Iteration**

loop

    Evaluate Policy

    Improve Policy

Converges because
policy always improves
and there are a finite
number of policies

**Value Iteration**

loop

    $V \leftarrow B[V]$

Converges because $B$ is
a contraction mapping

# Online MDP Planning

# Online MDP Planning

**Monte Carlo Tree Search**

# Online MDP Planning

**Monte Carlo Tree Search**

Search

Expand

Rollout

Backup

# Online MDP Planning

**Monte Carlo Tree Search**

Search

Expand

Rollout

Backup

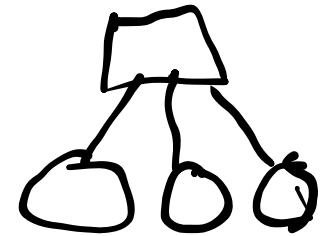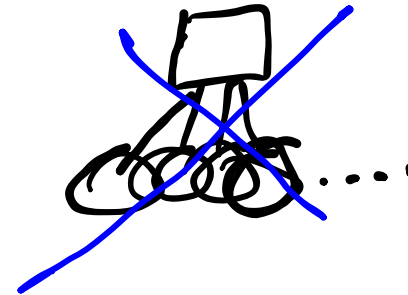$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}}$$

$Q(s,a)$

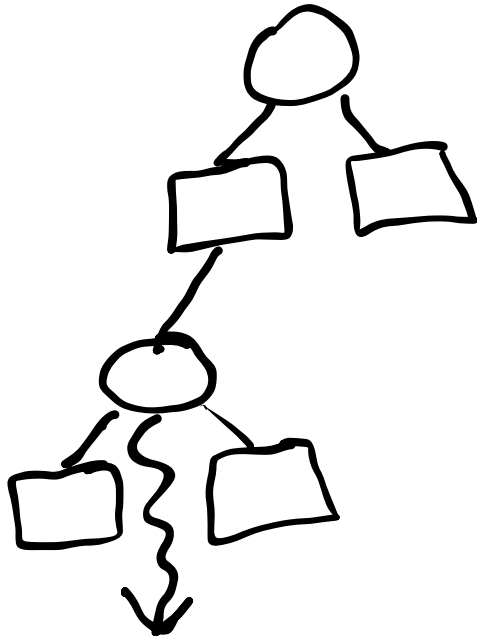# Online MDP Planning

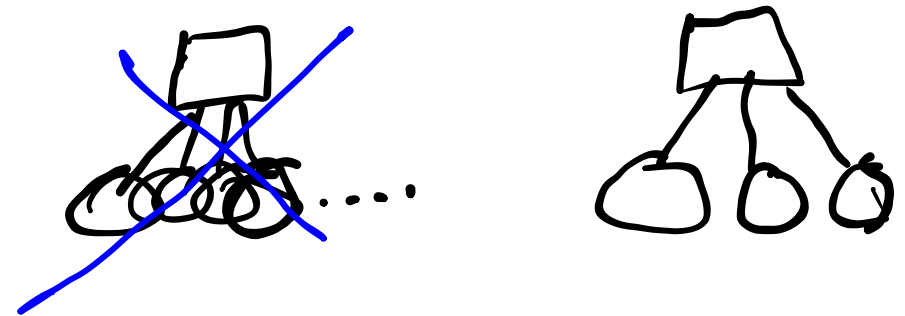**Monte Carlo Tree Search**                    **Sparse Sampling**

Search
Expand
Rollout
Backup

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}}$$

# Online MDP Planning

**Monte Carlo Tree Search**

**Sparse Sampling**

Search
Expand
Rollout
Backup

$$Q(s,a) + c \sqrt{\frac{\log N(s)}{N(s,a)}}$$

# Online MDP Planning

**Monte Carlo Tree Search**

**Sparse Sampling**

Search
Expand
Rollout
Backup

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}}$$

Guarantees *independent* of $|S|$!!

$s \in \mathbb{R}^n$
$a \in \mathbb{R}^m$

$\mathbb{R} \times \mathbb{R} \cdots$

# LQR

$$\mathbf{s}' = \mathbf{T}_s \mathbf{s} + \underline{\mathbf{T}_a} \mathbf{a} + \mathbf{w} \quad \overset{\text{Gaussian}}{\nwarrow}$$

$$R(\mathbf{s}, \mathbf{a}) = \underbrace{\mathbf{s}^\top \mathbf{R}_s \mathbf{s}} + \mathbf{a}^\top \mathbf{R}_a \mathbf{a}$$

$$\pi_h(\mathbf{s}) = -\underbrace{\left(\mathbf{T}_a^\top \mathbf{V}_{h-1} \mathbf{T}_a + \mathbf{R}_a\right)^{-1} \mathbf{T}_a^\top \mathbf{V}_{h-1} \mathbf{T}_s}_{K} \mathbf{s}$$

$U_h(\mathbf{s}) = s^\top V_h s + q_h$

$$\mathbf{V}_{h+1} = \underline{\mathbf{R}_s + \mathbf{T}_s^\top \mathbf{V}_h^\top \mathbf{T}_s - \left(\mathbf{T}_a^\top \mathbf{V}_h \mathbf{T}_s\right)^\top \left(\mathbf{R}_a + \mathbf{T}_a^\top \mathbf{V}_h \mathbf{T}_a\right)^{-1} \left(\mathbf{T}_a^\top \mathbf{V}_h \mathbf{T}_s\right)}$$

$$V_1 = R_s$$

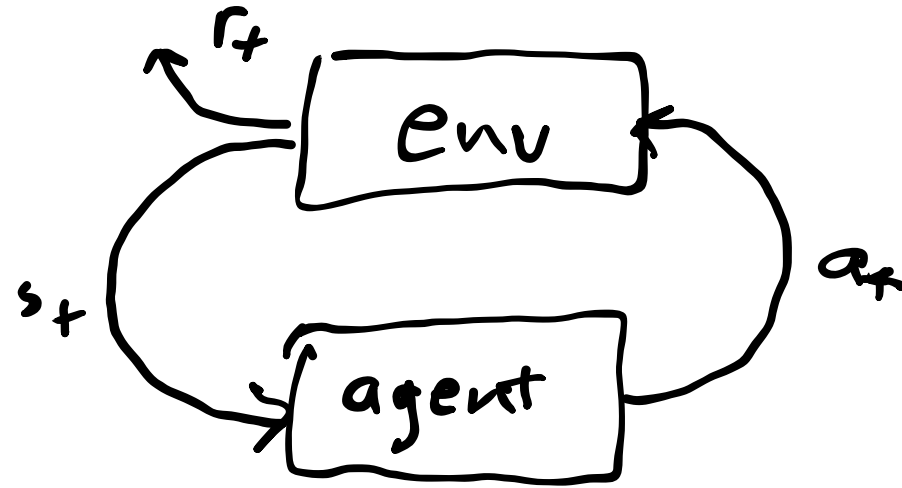# Reinforcement Learning

Challenges:

# Reinforcement Learning
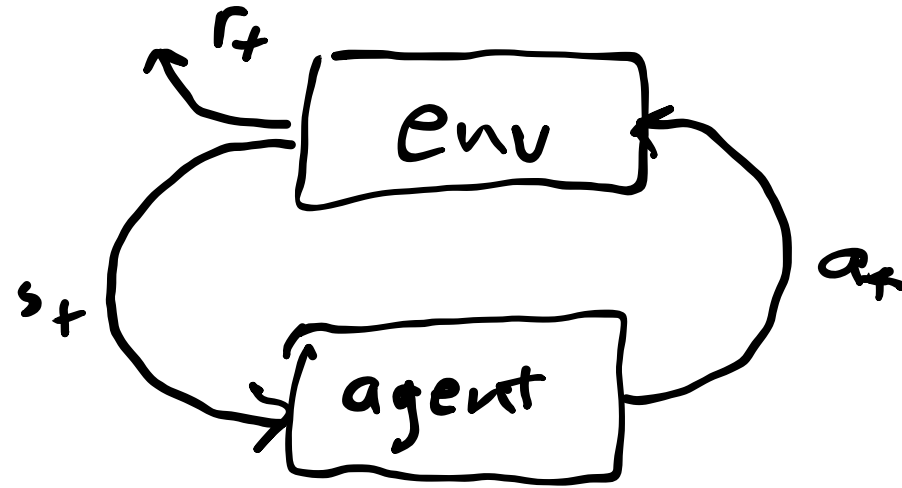


Challenges:

# Reinforcement Learning



Challenges:

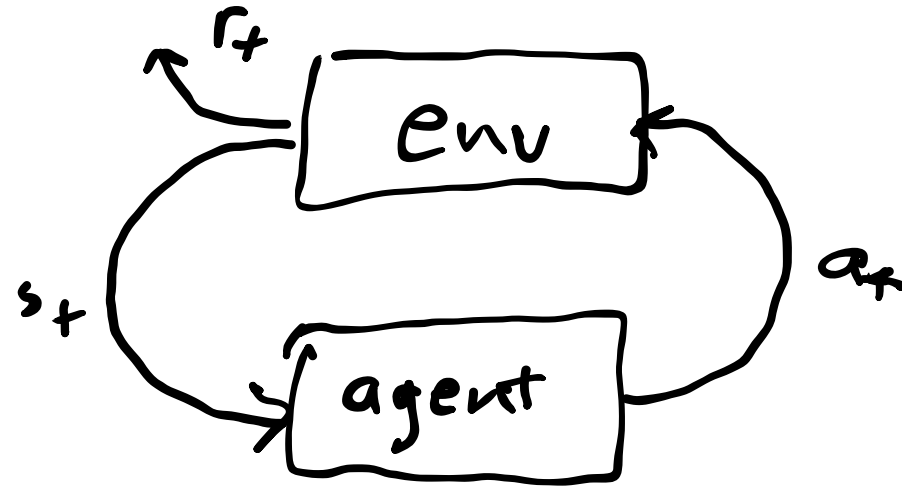    1. Exploration and Exploitation

# Reinforcement Learning



Challenges:

1. Exploration and Exploitation
2. Credit Assignment
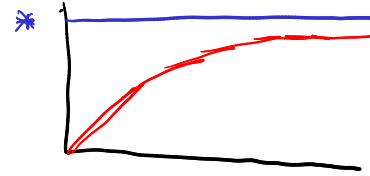
# Reinforcement Learning



Challenges:

    1. Exploration and Exploitation

    2. Credit Assignment

    3. Generalization

# Exploration
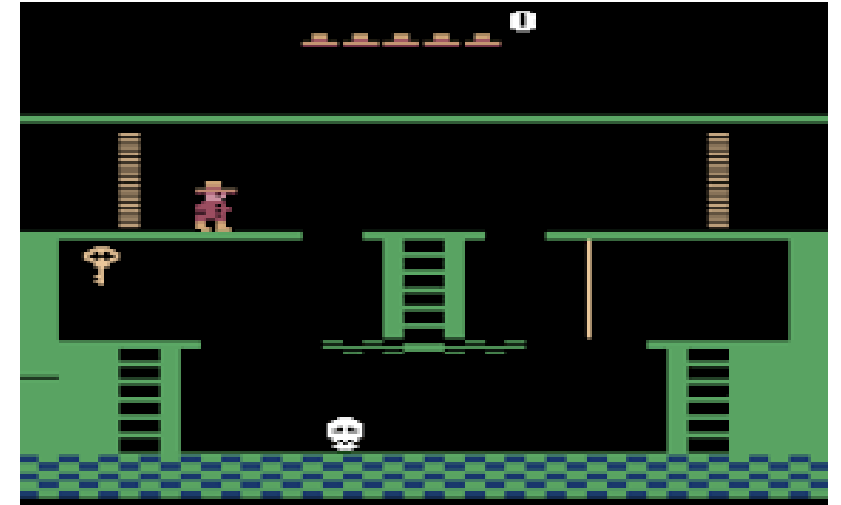
# Exploration

**Bandits**

- $\epsilon$-greedy
- softmax
- UCB
- Thompson Sampling
- Optimal DP Solution (solving a POMDP!)

logarithmic regret

# Exploration

**Bandits**

- $\epsilon$-greedy
- softmax
- UCB
- Thompson Sampling
- Optimal DP Solution (solving a POMDP!)
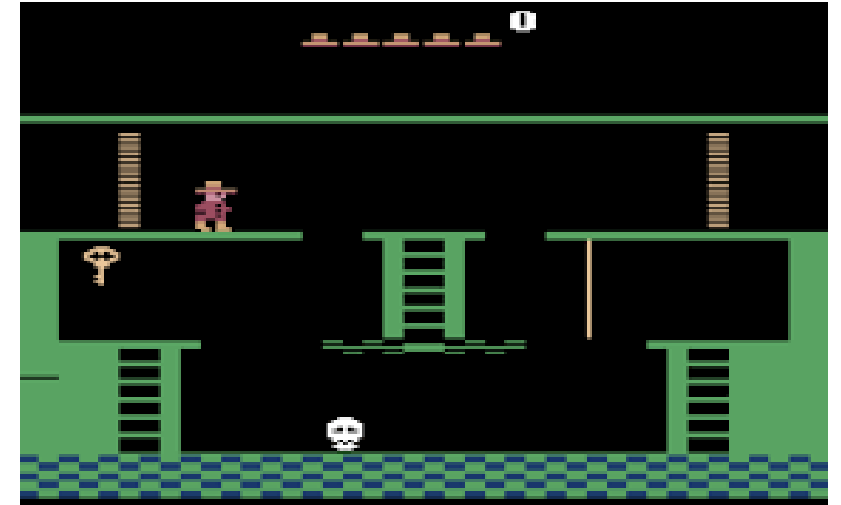


Montezuma's Revenge!

# Exploration

## Bandits

- $\epsilon$-greedy
- softmax
- UCB
- Thompson Sampling
- Optimal DP Solution (solving a POMDP!)
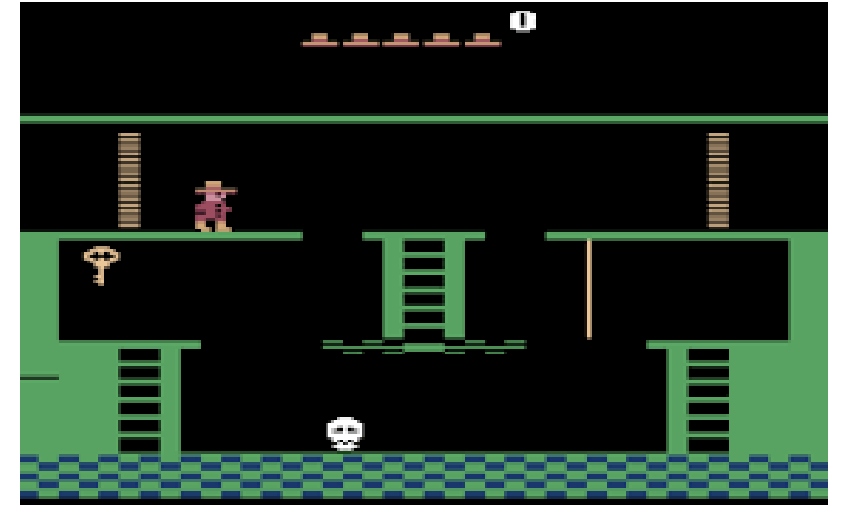
$$Q(s,a) + B(s,a)$$

- Pseudocounts



Montezuma's Revenge!

# Exploration

**Bandits**

- $\epsilon$-greedy
- softmax
- UCB
- Thompson Sampling
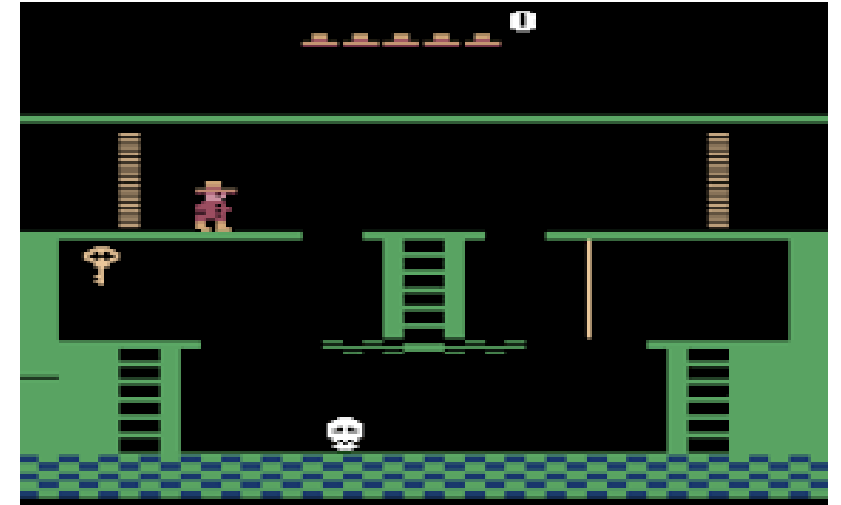- Optimal DP Solution (solving a POMDP!)



Montezuma's Revenge!

- Pseudocounts
- Curiosity: extra reward for bad prediction

# Exploration

**Bandits**

- $\epsilon$-greedy
- softmax
- UCB
- Thompson Sampling
- Optimal DP Solution (solving a POMDP!)



Montezuma's Revenge!

- Pseudocounts
- Curiosity: extra reward for bad prediction
- Random network distillation
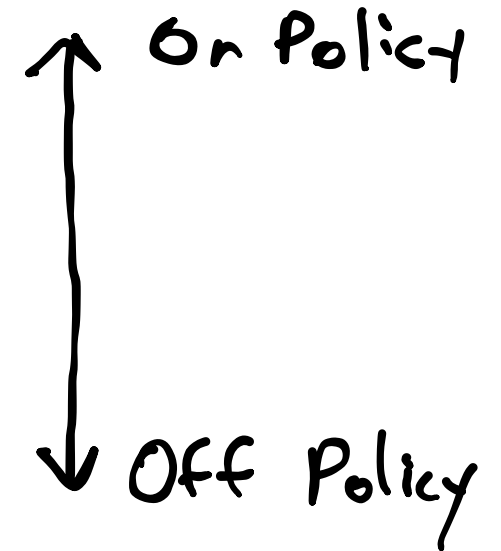
# RL Algorithms

# RL Algorithms

Model
Based

Model
Free

$\longleftrightarrow$

# RL Algorithms

Model
Based

Model
Free

Model Based ←————————————————→ Model Free

On Policy

Off Policy

# RL Algorithms

Model
Based

Model
Free

←————————————————————→

On Policy

↑
|
|
|
↓

Off Policy

MLMBTRL
(learn T,R)

# RL Algorithms

Model
Based

Model
Free

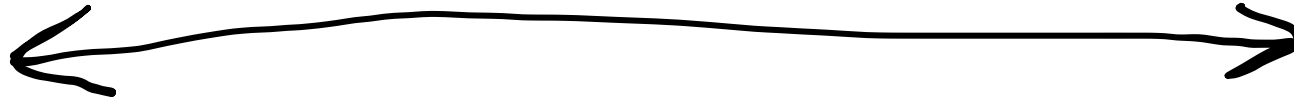learn Q    learn π    On Policy

MLMBTRL
(learn T,R)

Off Policy
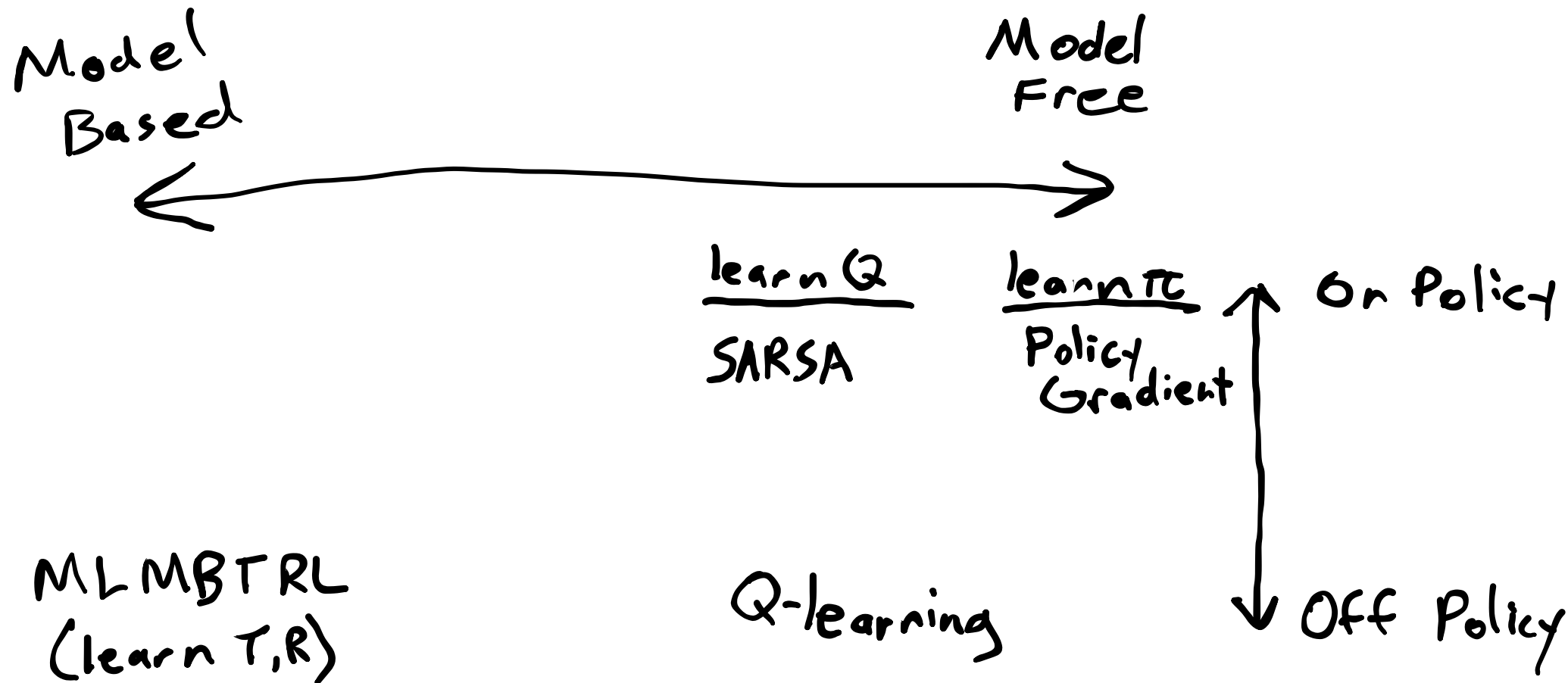
# RL Algorithms

Model
Based

Model
Free

learn Q

learn π
Policy
Gradient

On Policy

Off Policy

MLMBTRL
(learn T,R)

# RL Algorithms

Model
Based

Model
Free

$$\longleftrightarrow$$

$$\underline{\text{learn } Q}$$
SARSA

$$\underline{\text{learn } \pi}$$
Policy
Gradient

On Policy

Off Policy

MLMBTRL
(learn T,R)

# RL Algorithms

Model
Based

Model
Free



$$\frac{\text{learn } Q}{\text{SARSA}}$$

$$\frac{\text{learn } \pi}{\text{Policy Gradient}}$$

On Policy

Off Policy

MLMBTRL
(learn T, R)

Q-learning

# RL Algorithms



Model Based ⟷ Model Free

learn Q
―――
SARSA

learn π
Policy Gradient

On Policy
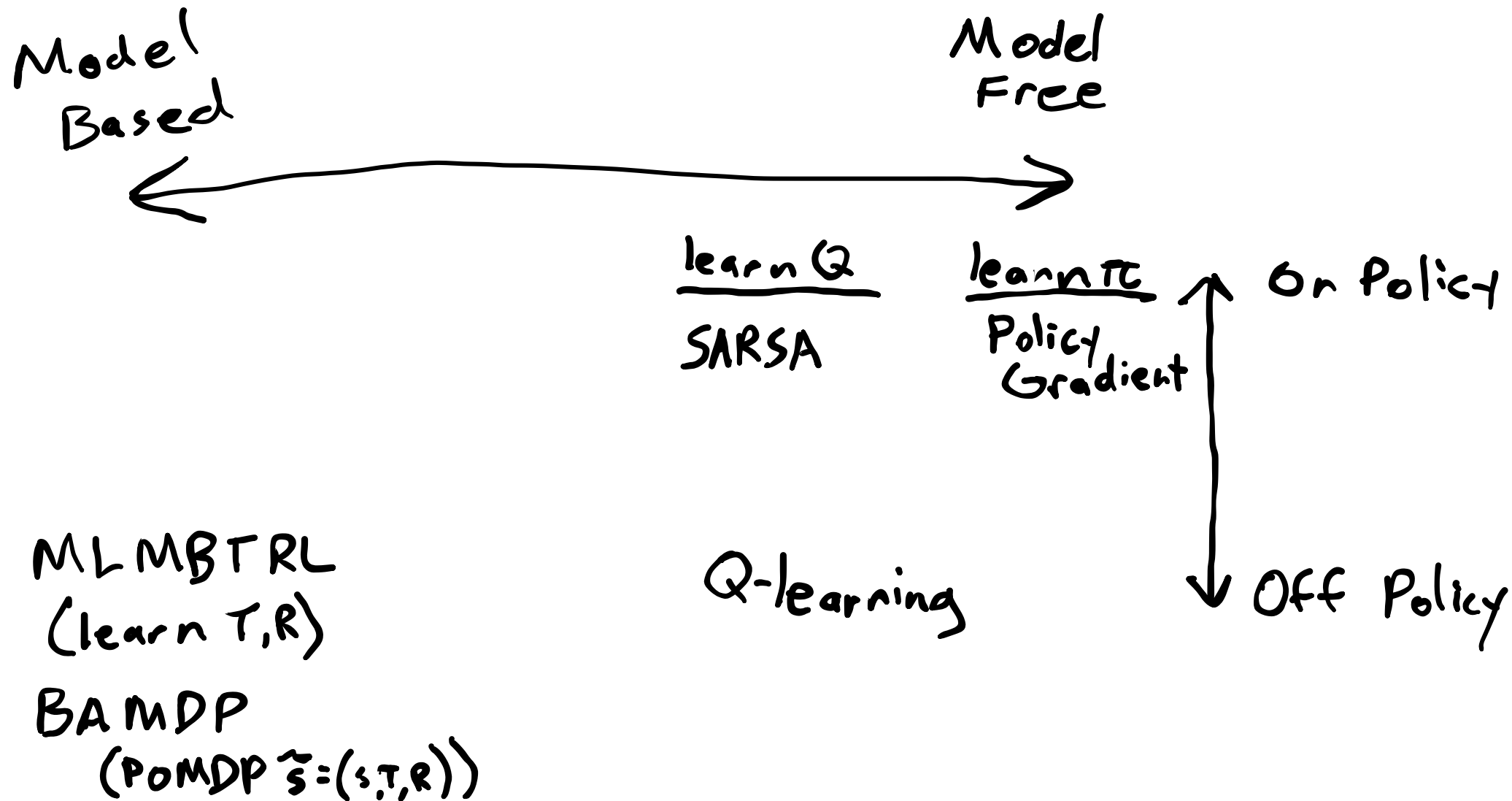
Off Policy

MLMBTRL
(learn T,R)

BAMDP
(POMDP $\tilde{s} = (s, T, R)$)

Q-learning

# Policy Gradient

# Policy Gradient

- Likelihood ratio trick

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

# Policy Gradient

- Likelihood ratio trick
- Causality

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^{(k)} \mid s^{(k)}) \gamma^{k-1} \left( r_{\text{to-go}}^{(k)} - r_{\text{base}}(s^{(k)}) \right) \right]$$
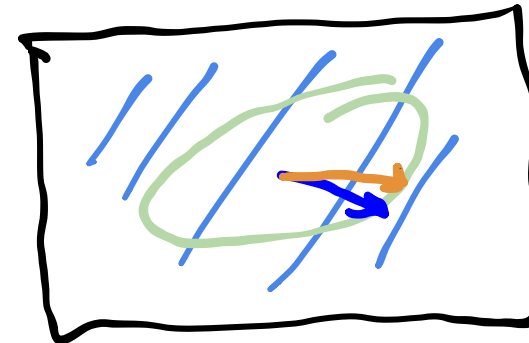
# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

$$\nabla_\theta \, p_\theta(\tau) \; = p_\theta(\tau) \, \nabla_\theta \, \log p_\theta(\tau)$$

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^{(k)} \mid s^{(k)}) \gamma^{k-1} \left( r_{\text{to-go}}^{(k)} - r_{\text{base}}(s^{(k)}) \right) \right]$$

- Natural Gradient

# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^{(k)} \mid s^{(k)}) \gamma^{k-1} \left( r_{\text{to-go}}^{(k)} - r_{\text{base}}(s^{(k)}) \right) \right]$$

- Natural Gradient

KL div. Bound

# Q-Learning

# Q-Learning

**SARSA**

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r_t + \gamma Q(s', a') - Q(s,a))$$

# Q-Learning

**SARSA**

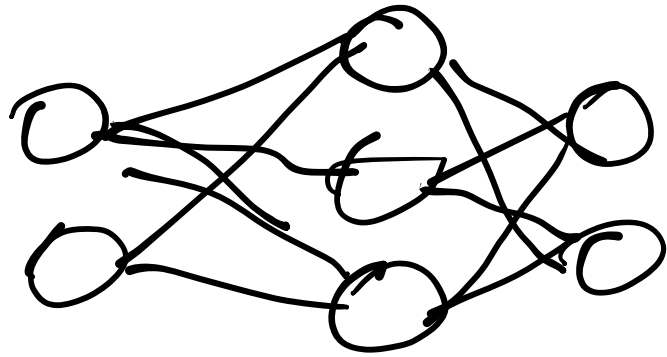$$Q(s, a) \leftarrow Q(s, a) + \alpha(r_t + \gamma Q(s', a') - Q(s, a))$$

Eligibility Traces

# Q-Learning

**SARSA**

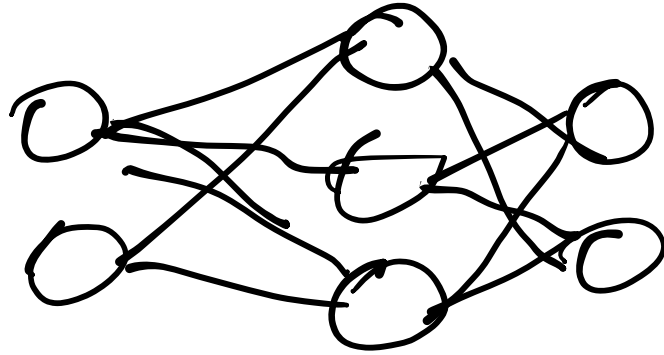$$Q(s,a) \leftarrow Q(s,a) + \alpha(r_t + \gamma Q(s',a') - Q(s,a))$$

Eligibility Traces

**Q-learning**

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

# Q-Learning

**SARSA**

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r_t + \gamma Q(s', a') - Q(s, a))$$

Eligibility Traces

**Q-learning**

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

Double Q Learning

# Neural Networks and DQN

# Neural Networks and DQN

# Neural Networks and DQN



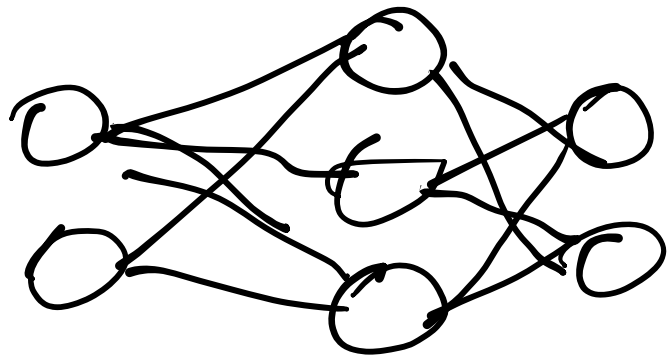$$f_\theta(x) = \sigma(W_2\sigma(W_1x + b_1) + b_2)$$

# Neural Networks and DQN
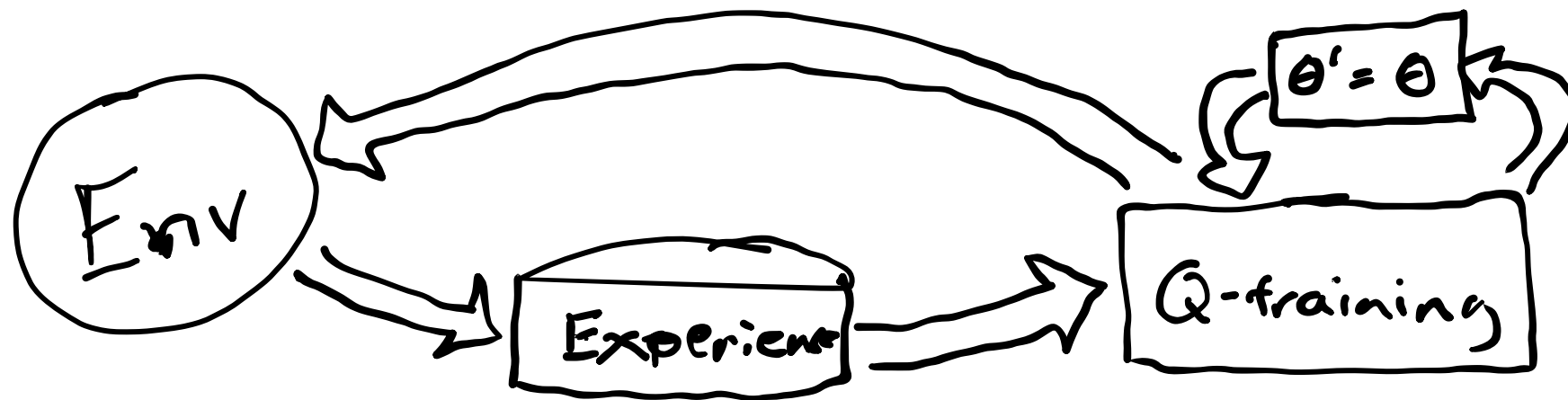
$$f_\theta(x) = \sigma(W_2\sigma(W_1x + b_1) + b_2)$$

Backprop

# Neural Networks and DQN

$$f_\theta(x) = \sigma(W_2 \sigma(W_1 x + b_1) + b_2)$$

Backprop

# Actor-Critic

# Actor-Critic

- Actor: $\pi_\theta$

# Actor-Critic

- Actor: $\pi_\theta$
- Critic: $Q_\phi$

# Actor-Critic

- Actor: $\pi_\theta$
- Critic: $Q_\phi$

**Soft Actor Critic**

# Actor-Critic

- Actor: $\pi_\theta$
- Critic: $Q_\phi$

**Soft Actor Critic**

$$J(\pi) = E\left[\sum_{t=0}^{\infty} \gamma^t \left(r_t + \alpha \mathcal{H}(\pi(\cdot \mid s_t)))\right)\right]$$

# POMDPs

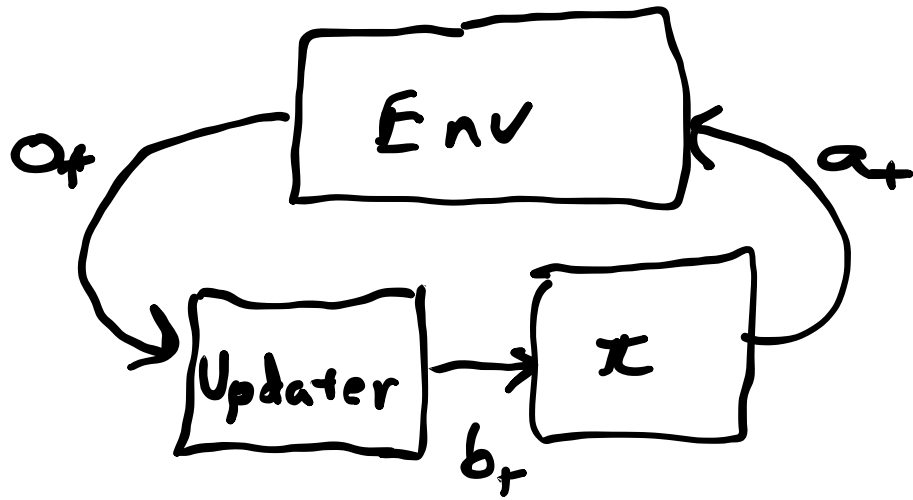# POMDPs

$$(S, A, T, R, O, Z, \gamma)$$

# POMDPs

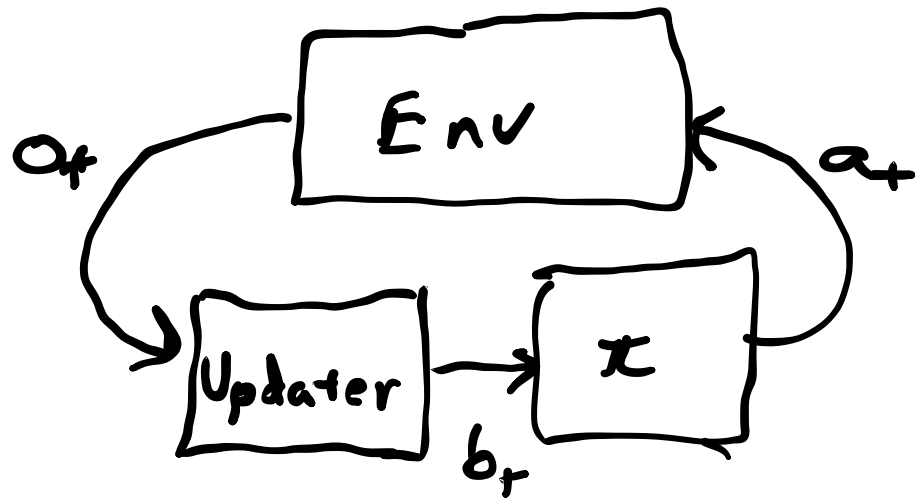$$(S, A, T, R, O, Z, \gamma)$$

# POMDPs

$$(S, A, T, R, O, Z, \gamma)$$



**Belief Updates**

- Discrete Bayesian Filter
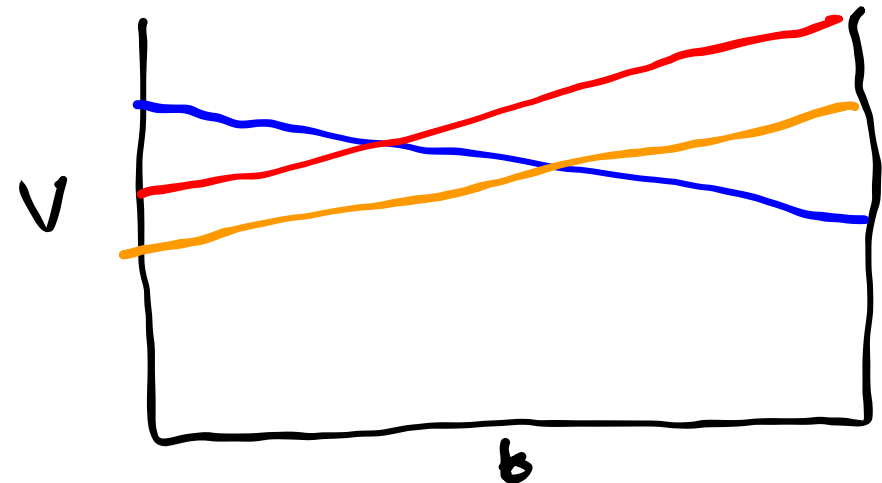- Particle Filter

# POMDPs
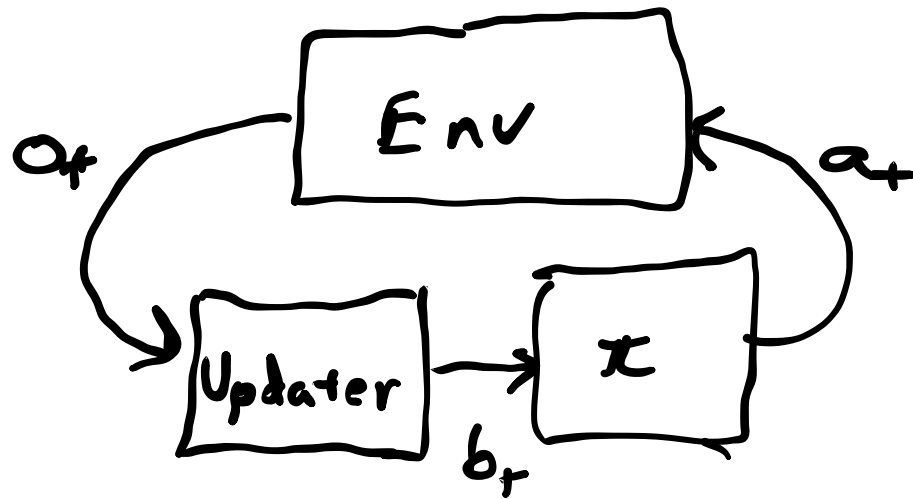
$$(S, A, T, R, O, Z, \gamma)$$



**Belief Updates**

- Discrete Bayesian Filter
- Particle Filter

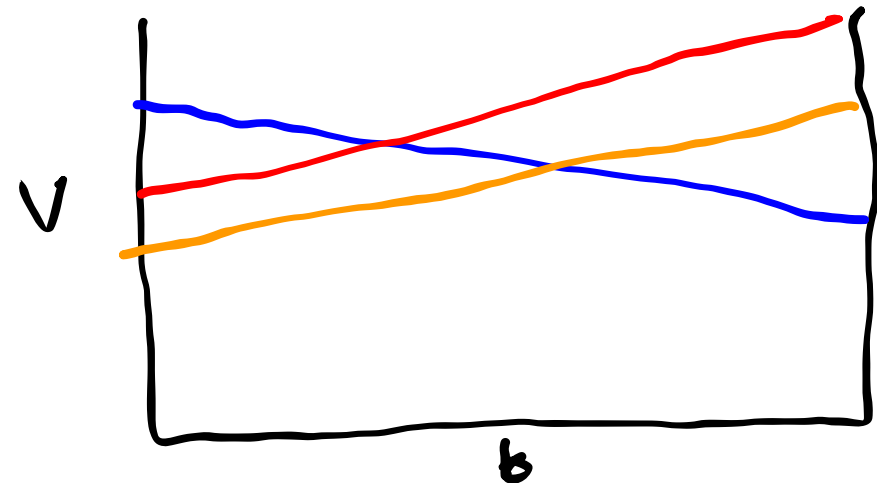**Alpha Vectors**

# POMDPs

$$(S, A, T, R, O, Z, \gamma)$$



**Belief Updates**

- Discrete Bayesian Filter
- Particle Filter

**Alpha Vectors**

- Each alpha vector corresponds to a conditional plan
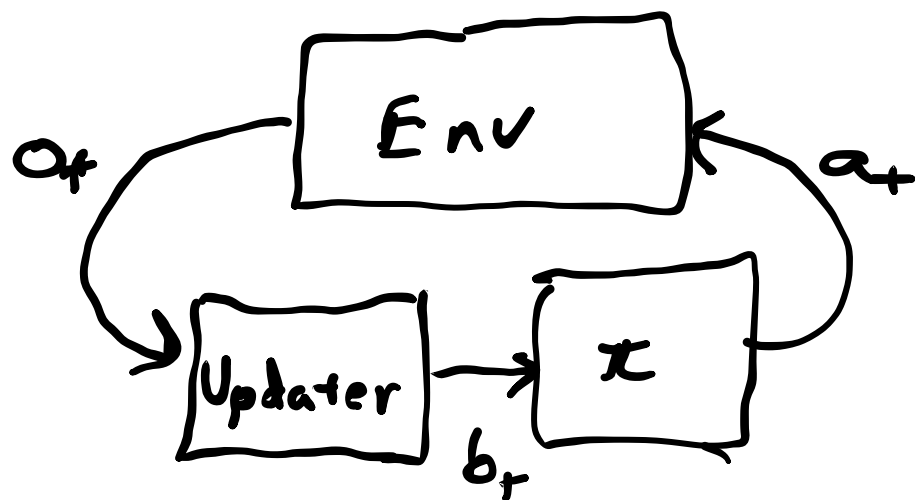
# POMDPs

$$(S, A, T, R, O, Z, \gamma)$$



**Belief Updates**

- Discrete Bayesian Filter
- Particle Filter

**Alpha Vectors**



- Each alpha vector corresponds to a conditional plan
- You can prune alpha vectors by solving an LP

# POMDP Approximations

# POMDP Approximations

**Formulation**

- Certainty Equivalence $\longleftarrow$ Optimal for LQG
- QMDP

$$\pi_{QMDP}(s) = \arg\max_a \mathbb{E}_{s \sim b}\left[Q_{MDP}(s,a)\right]$$

# POMDP Approximations

**Formulation**

- Certainty Equivalence
- QMDP

**Numerical**
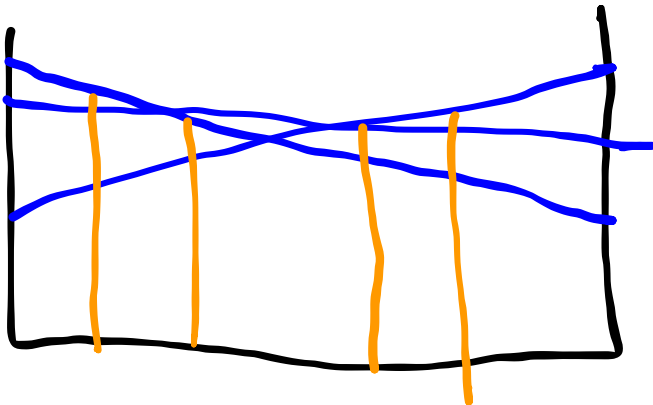
# POMDP Approximations

## Formulation

- Certainty Equivalence
- QMDP

## Numerical

### Offline

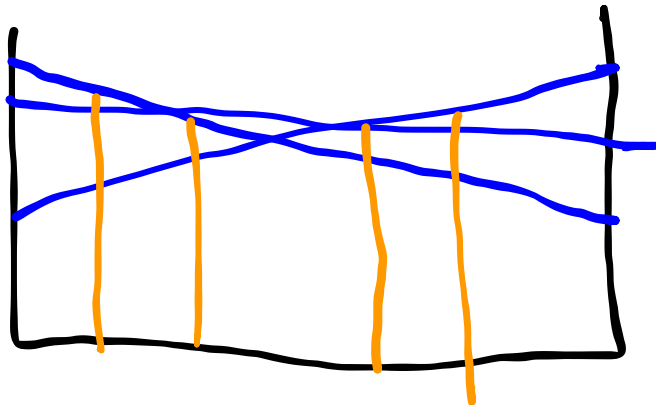- Point-Based Value Iteration
- SARSOP

# POMDP Approximations

## Formulation

- Certainty Equivalence
- QMDP

## Numerical

### Offline

- Point-Based Value Iteration
- SARSOP

### Online

- POMCP
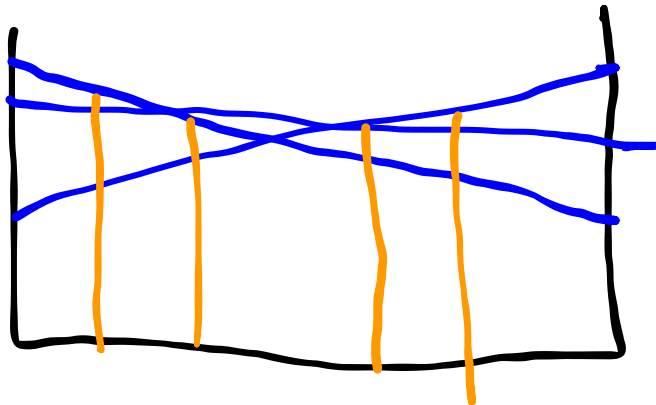- DESPOT

# POMDP Approximations

## Formulation

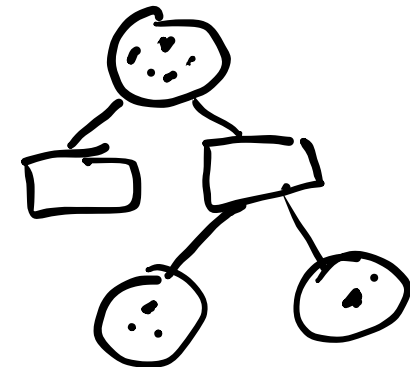- Certainty Equivalence
- QMDP

## Numerical

### Offline

- Point-Based Value Iteration
- SARSOP



### Online

- POMCP
- DESPOT

# Simple Games

# Simple Games

- Optimal Solutions

# Simple Games

- ~~Optimal Solutions~~ No!

# Simple Games

- ~~Optimal Solutions~~ No!
- Equilibria (e.g. Nash Equilibria)

# Simple Games

- ~~Optimal Solutions~~ No!
- Equilibria (e.g. Nash Equilibria)

# Simple Games

- ~~Optimal Solutions~~ **No!**
- Equilibria (e.g. Nash Equilibria)

|        |       |
|--------|-------|
| -1,-1  | -3,0  |
| 0,-3   | -2,-2 |

- Every finite game has at least 1 Nash Equilibrium

# Simple Games

- ~~Optimal Solutions~~ **No!**
- Equilibria (e.g. Nash Equilibria)

| | |
|---|---|
| -1,-1 | -3,0 |
| 0,-3 | -2,-2 |

- Every finite game has at least 1 Nash Equilibrium
- Might be pure or mixed

# Simple Games

- ~~Optimal Solutions~~ **No!**
- Equilibria (e.g. Nash Equilibria)
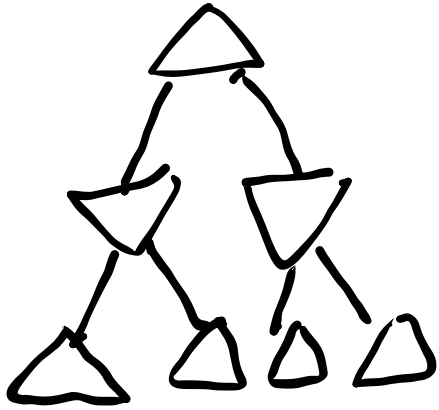


| | |
|---|---|
| -1,-1 | -3,0 |
| 0,-3 | -2,-2 |

- Every finite game has at least 1 Nash Equilibrium
- Might be pure or mixed
- Algorithms like fictitious play converge in special cases
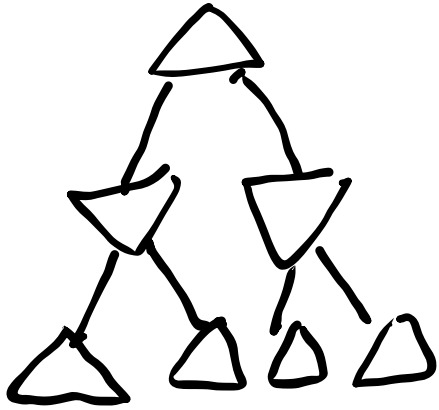
# Turn Taking Games

# Turn Taking Games
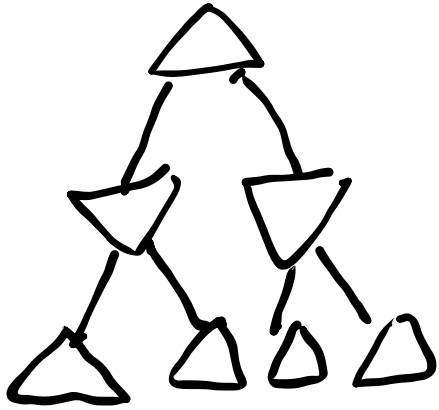
# Turn Taking Games



- Value Function Backup
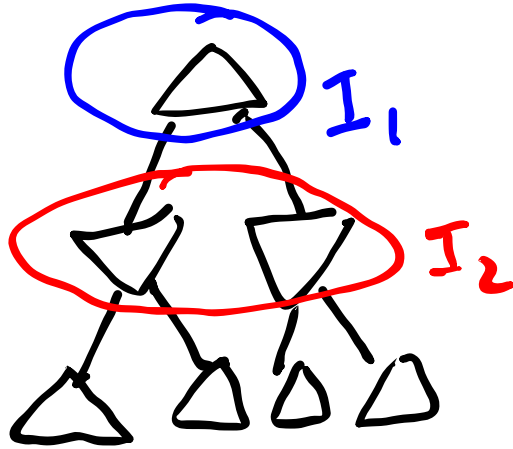
# Turn Taking Games

- Value Function Backup
- $\alpha\beta$ Pruning

# Turn Taking Games

- Value Function Backup
- $\alpha\beta$ Pruning
- Incomplete Information Extensive Form

# Turn Taking Games



- Value Function Backup
- $\alpha\beta$ Pruning
- Incomplete Information Extensive Form

# Markov Games and POMGS

# Markov Games and POMGS

**Markov Games**

- All players play simultaneously
- Transitions are stochastic
- Best response involves solving an MDP
- Can be reduced to a simple game with policies as actions
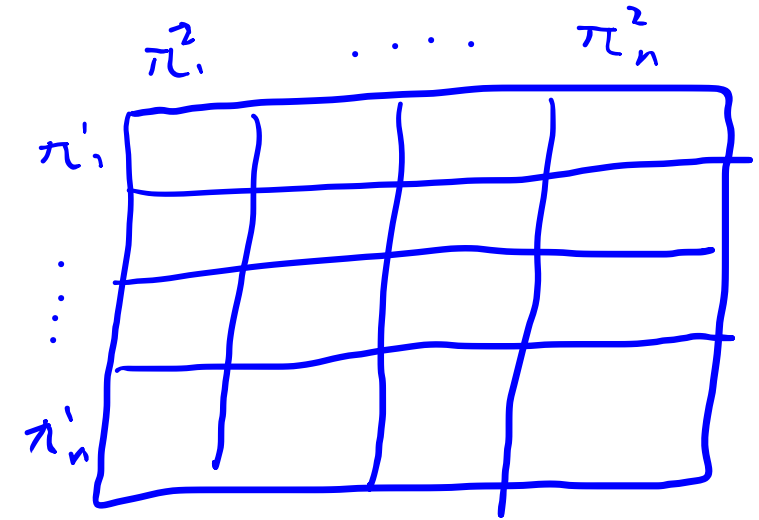
# Markov Games and POMGS

**Markov Games**

- All players play simultaneously
- Transitions are stochastic
- Best response involves solving an MDP
- Can be reduced to a simple game with policies as actions

# Markov Games and POMGS
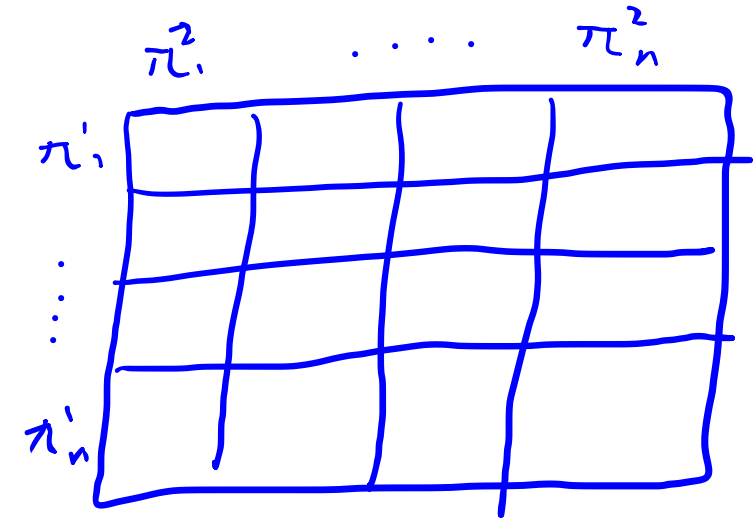
## Markov Games

- All players play simultaneously
- Transitions are stochastic
- Best response involves solving an MDP
- Can be reduced to a simple game with policies as actions

## Partially Observable Markov Games

- Each player receives a noisy observation at each step
- Beliefs not practical to compute
- Can be reduced to simple game with policies as actions

# Markov Games and POMGS

## Markov Games

- All players play simultaneously
- Transitions are stochastic
- Best response involves solving an MDP
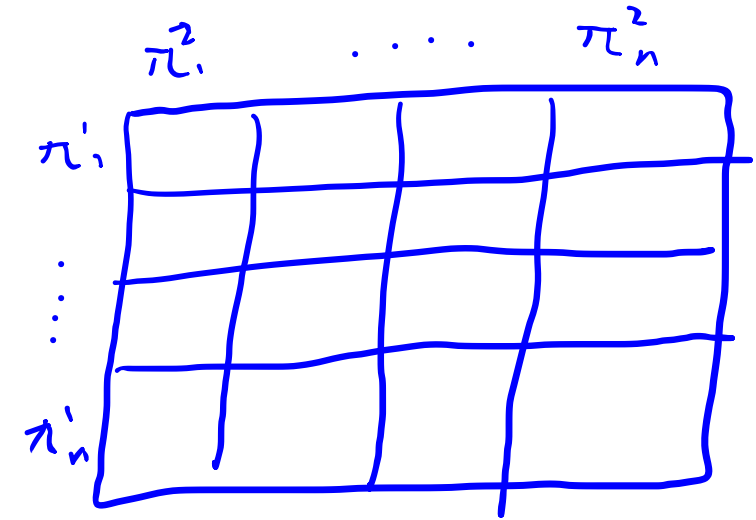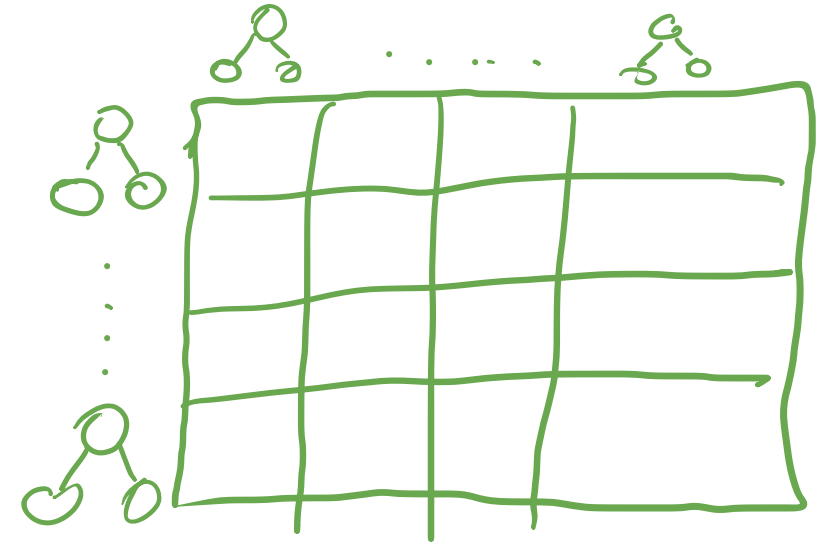- Can be reduced to a simple game with policies as actions

## Partially Observable Markov Games

- Each player receives a noisy observation at each step
- Beliefs not practical to compute
- Can be reduced to simple game with policies as actions

# Fictitious Play in Markov Games

# Recap

# Recap

**After DMU you have basic tools to deal with 4 Big Problems:**

# Recap

**After DMU you have basic tools to deal with 4 Big Problems:**

1. Immediate and Future Rewards

# Recap

**After DMU you have basic tools to deal with 4 Big Problems:**

1. Immediate and Future Rewards
2. Unknown Models

# Recap

**After DMU you have basic tools to deal with 4 Big Problems:**

1. Immediate and Future Rewards
2. Unknown Models
3. Partial Observability

# Recap

**After DMU you have basic tools to deal
with 4 Big Problems:**

1. Immediate and Future Rewards
2. Unknown Models
3. Partial Observability
4. Other Agents

ASEN 6519
DMU++