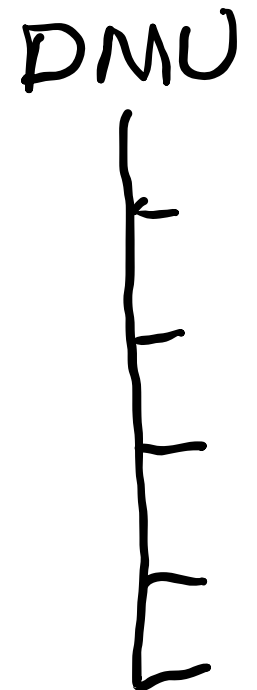


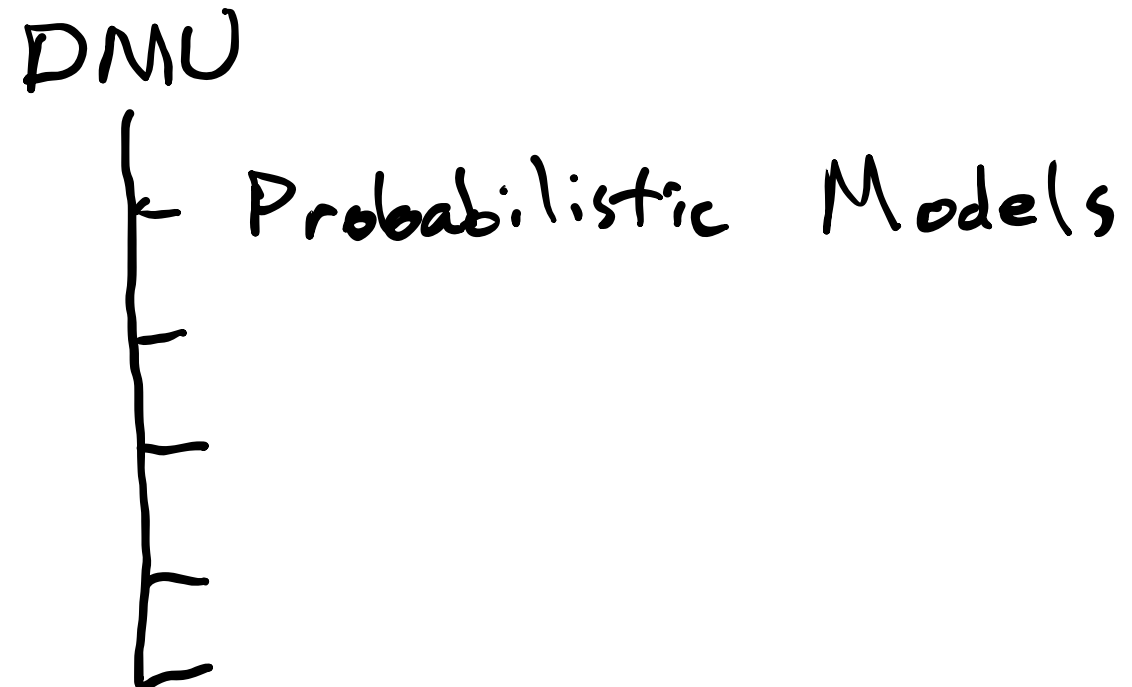
# Recap

# Recap



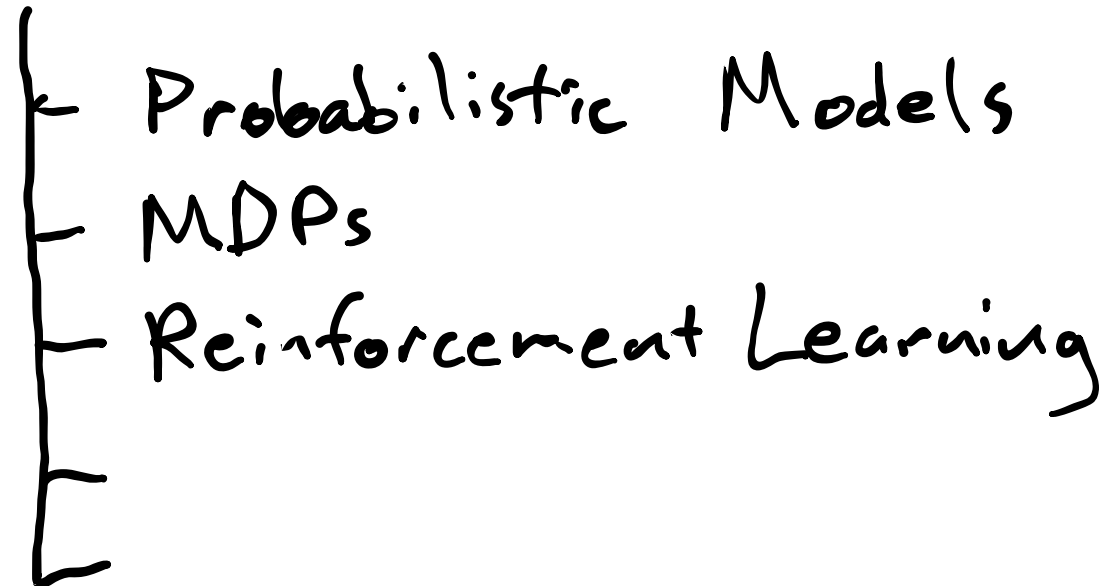
# Recap

DMU



# Recap

DMU



# Recap

DMU

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

h

# Probabilistic Models

# Probabilistic Models

3 Rules

$P(A)$   
 $P(A, B)$   
 $P(A|B)$

# Probabilistic Models

## 3 Rules

$$\begin{aligned} P(A) \\ P(A, B) \\ P(A|B) \end{aligned}$$

$$\begin{aligned} 1. \quad & 0 \leq P(X | Y) \leq 1 \\ & \sum_{x \in X} P(x | Y) = 1 \end{aligned}$$



# Probabilistic Models

$P(A)$   
 $P(A, B)$   
 $P(A|B)$

## 3 Rules

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

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

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

# Probabilistic Models

$$P(A)$$
$$P(A, B)$$
$$P(A|B)$$

## 3 Rules

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

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

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

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

# Probabilistic Models

$$\begin{aligned} P(A) \\ P(A, B) \\ P(A|B) \end{aligned}$$

## 3 Rules

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

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

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

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

## Bayes Rule

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

# Probabilistic Models

$$\begin{aligned} P(A) \\ P(A, B) \\ P(A|B) \end{aligned}$$

## 3 Rules

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

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

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

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

## Bayes Rule

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

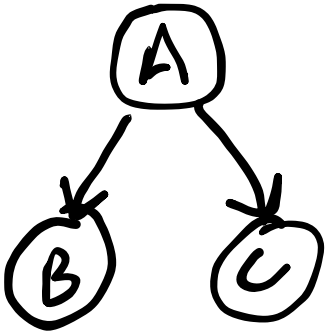
## Independence

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

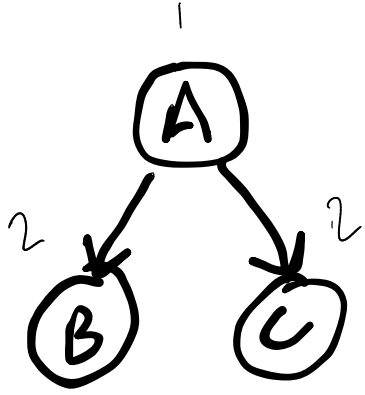
$$A \perp B | C \iff P(A, B | C) = P(A | C)P(B | C)$$

# Bayesian Networks

# Bayesian Networks



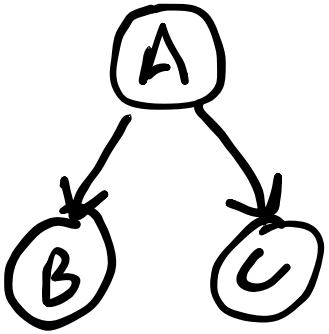
# Bayesian Networks



**Chain Rule**

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

# Bayesian Networks



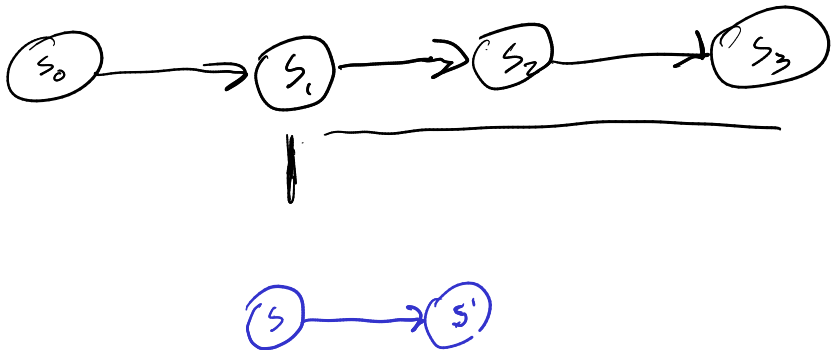
## Chain Rule

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

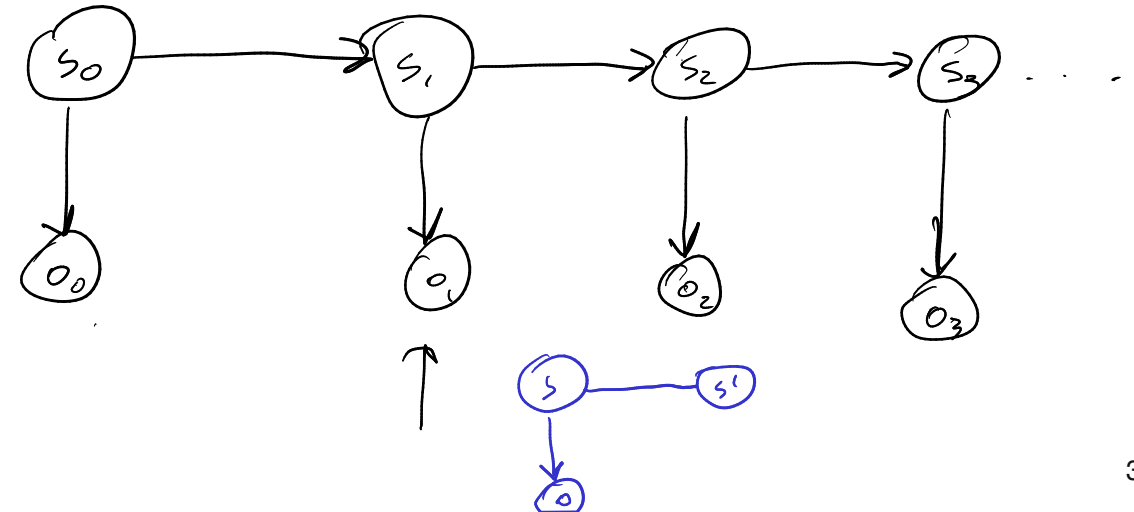


## Conditional Independence

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



$$o_2 \not\perp o_0 \mid o_1$$





# 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

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

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

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

# Markov Decision Processes

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

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

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

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

# Markov Decision Processes

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

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

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

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

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

# Markov Decision Processes

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

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

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

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

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

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

# Markov Decision Processes

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

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

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

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

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

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

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma E[V^{\pi}(s')] \quad \leftarrow \text{policy evaluation}$$

$$\rightarrow V^*(s) = \max_a \{ R(s, a) + \gamma E[V^*(s')] \} \quad \leftarrow \begin{array}{l} \text{certificate} \\ \text{of optimality} \\ \text{"Bellman's Eq."} \end{array}$$

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



# Markov Decision Processes

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

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

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


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

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

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

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

Policy Evaluation


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

Bellman's Equation: Certificate of Optimality

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

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

*Curse of Dimensionality  $\rightarrow$  Large  $|S|$*

# Online MDP Planning

# Online MDP Planning

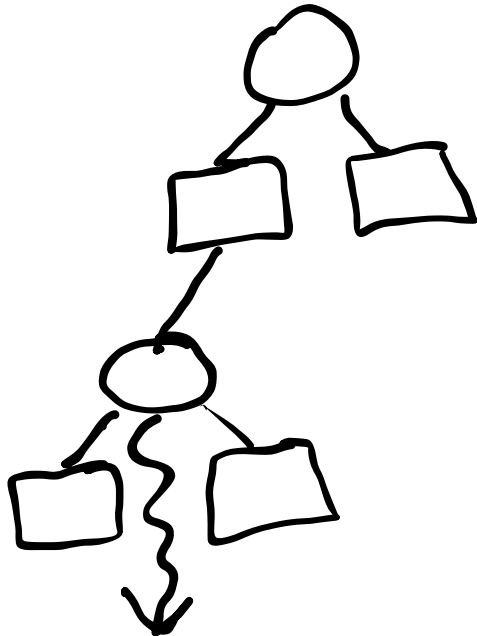
**Monte Carlo Tree Search**



# Online MDP Planning

## Monte Carlo Tree Search

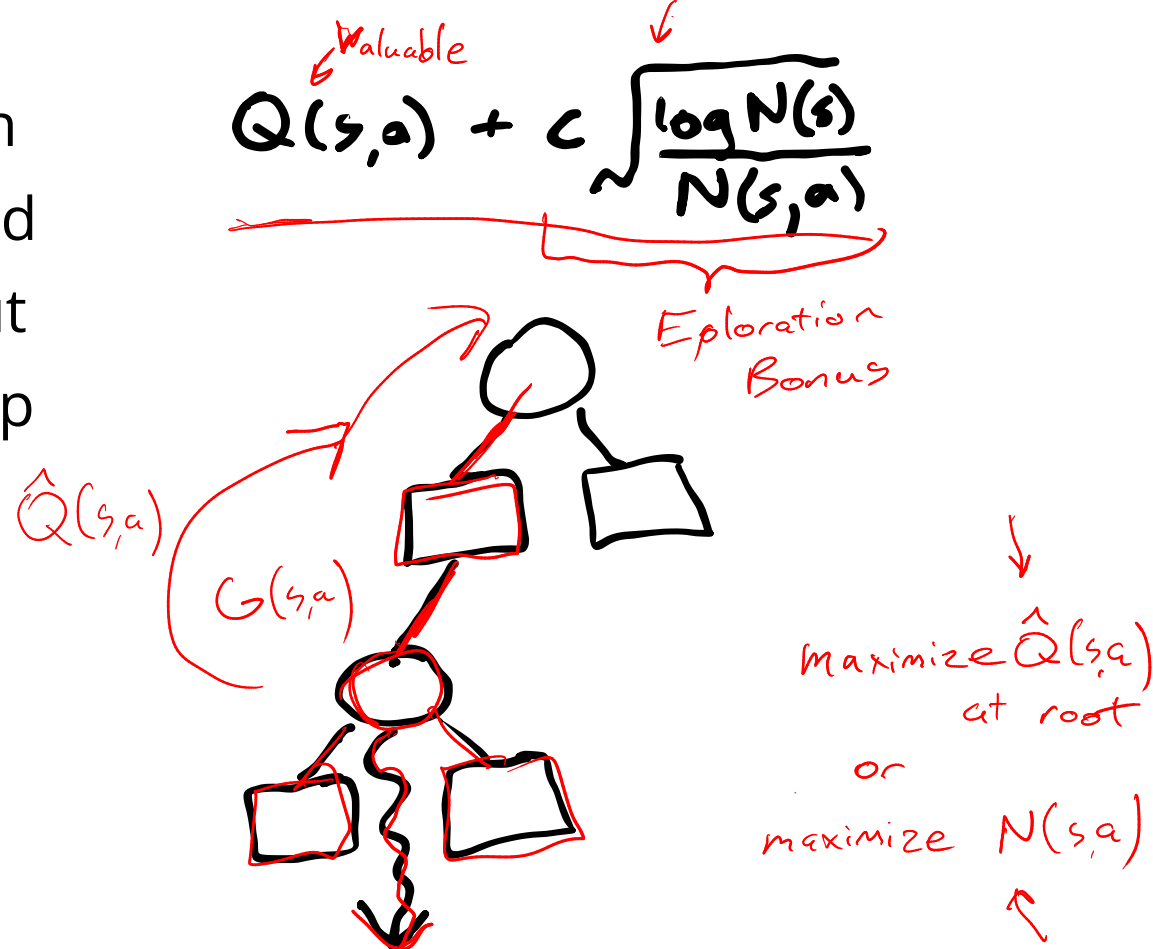
Search  
Expand  
Rollout  
Backup



# Online MDP Planning

## Monte Carlo Tree Search

Search  
Expand  
Rollout  
Backup

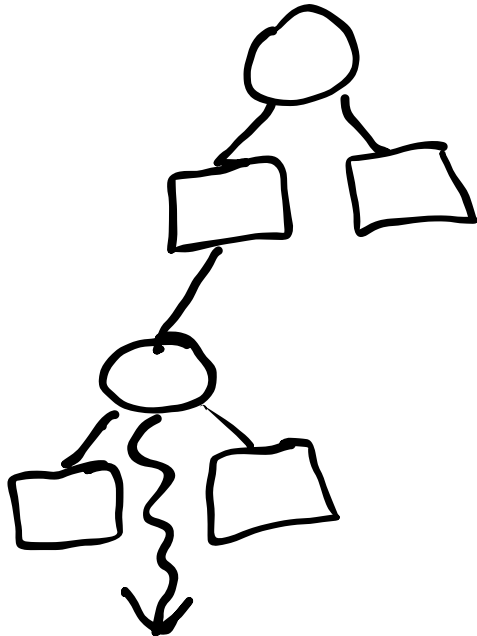


# Online MDP Planning

## Monte Carlo Tree Search

Search  
Expand  
Rollout  
Backup

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



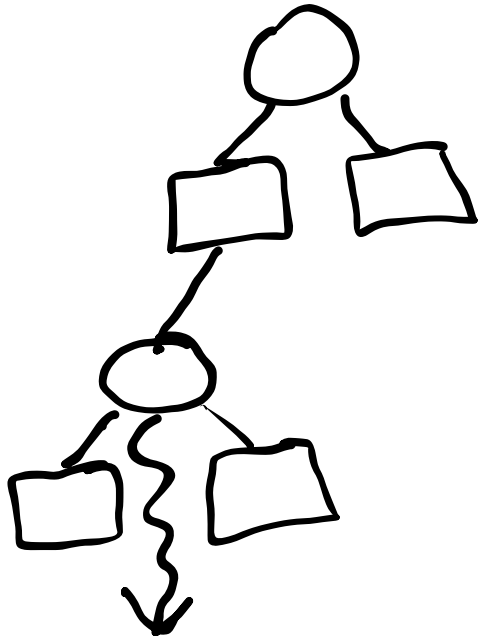
## Sparse Sampling

# Online MDP Planning

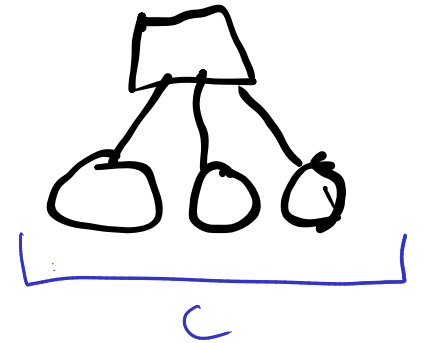
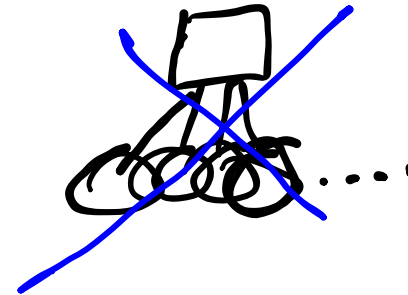
## Monte Carlo Tree Search

Search  
Expand  
Rollout  
Backup

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



## Sparse Sampling

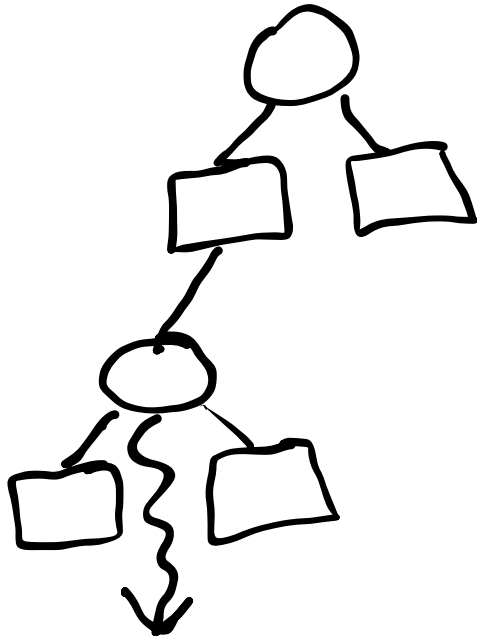


# Online MDP Planning

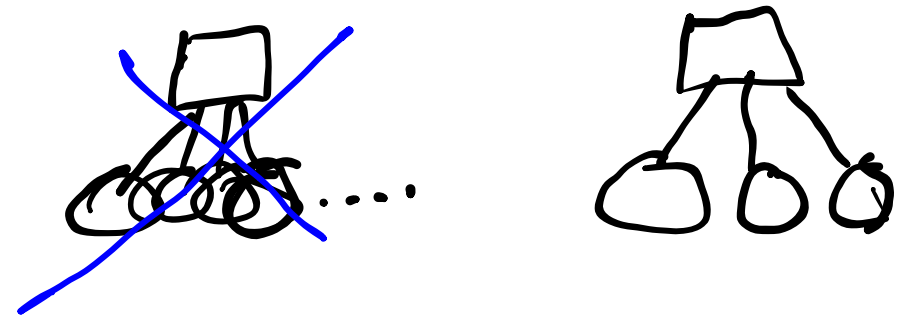
## Monte Carlo Tree Search

Search  
Expand  
Rollout  
Backup

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



## Sparse Sampling

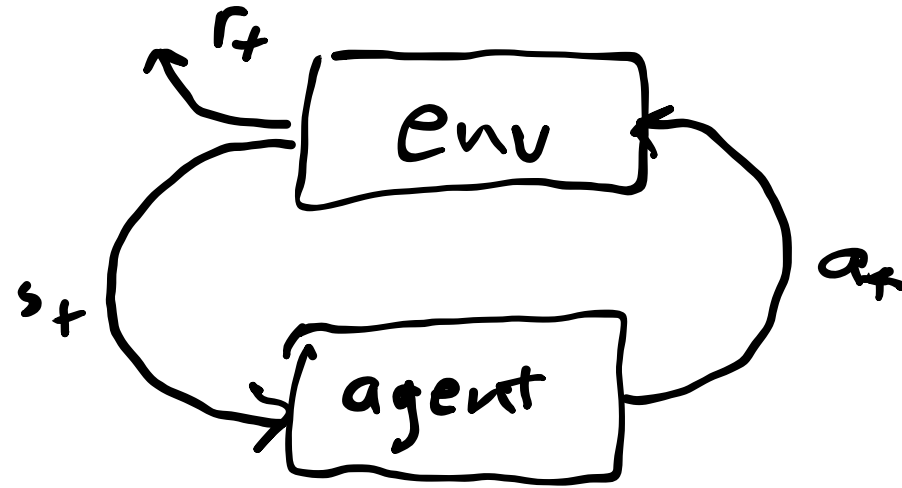


Guarantees *independent* of  $|S|!!$

# Reinforcement Learning

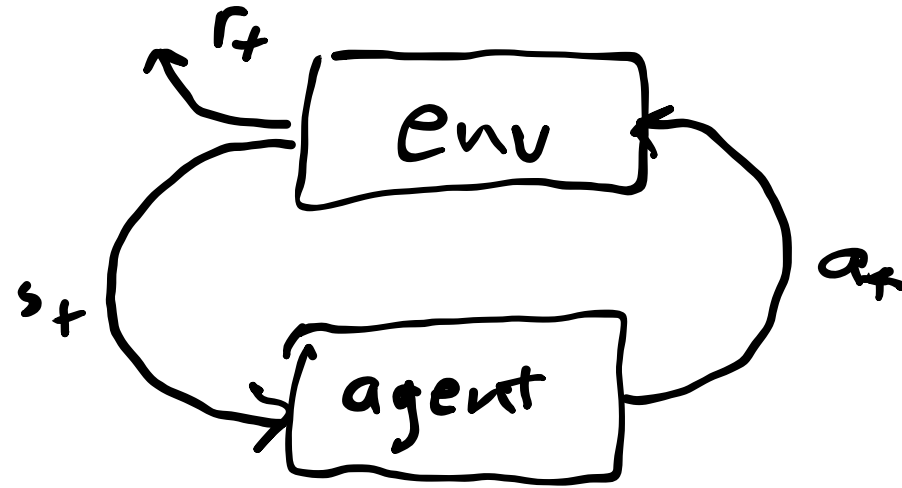
Challenges:

# Reinforcement Learning



Challenges:

# Reinforcement Learning

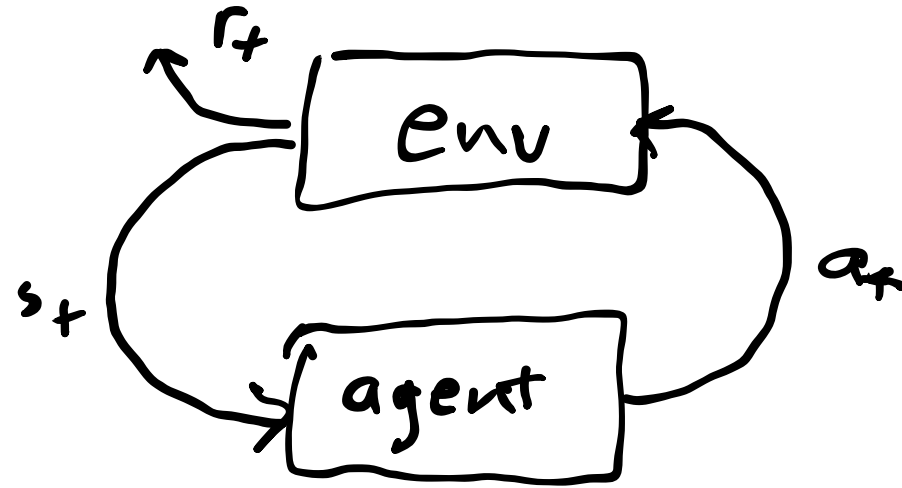


Challenges:

1. Exploration and Exploitation



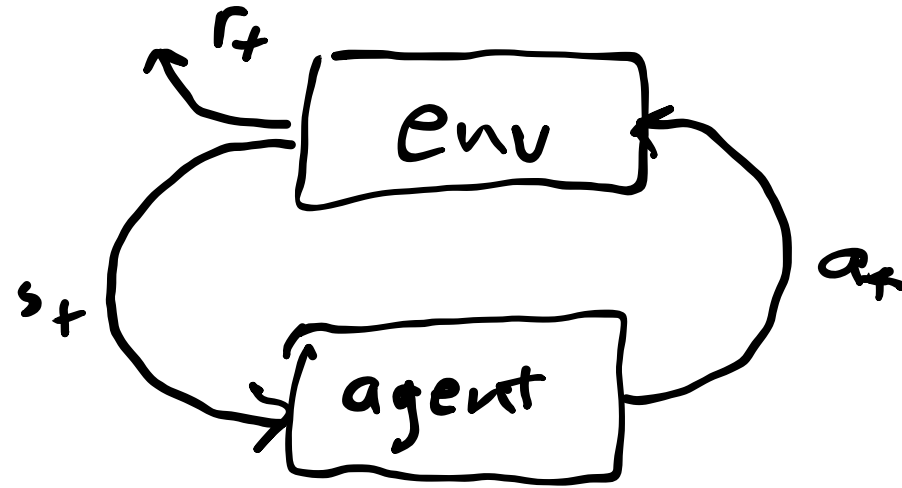
# Reinforcement Learning



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 *easiest*

- softmax

- UCB

- Thompson Sampling

- Optimal DP Solution (solving a POMDP!)

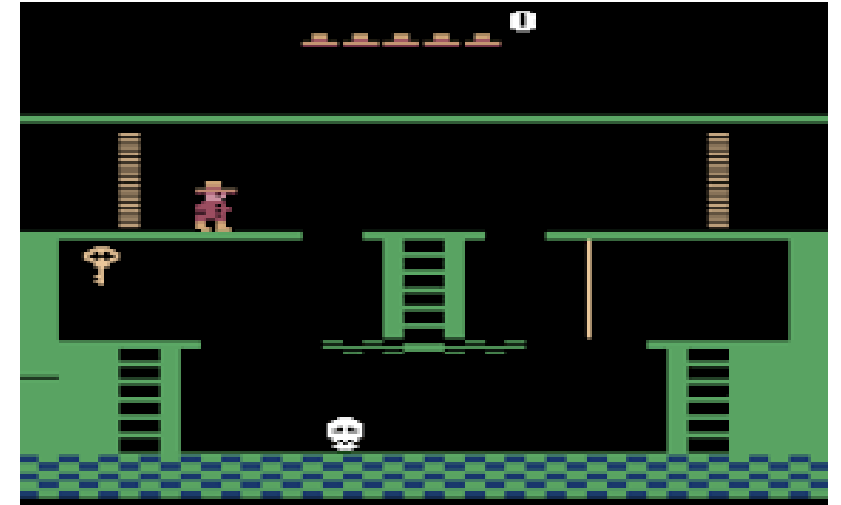
$$\hat{Q}(a) + c \sqrt{\frac{\log N}{N(a)}}$$

logarithmic  
regret

# Exploration

## Bandits

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



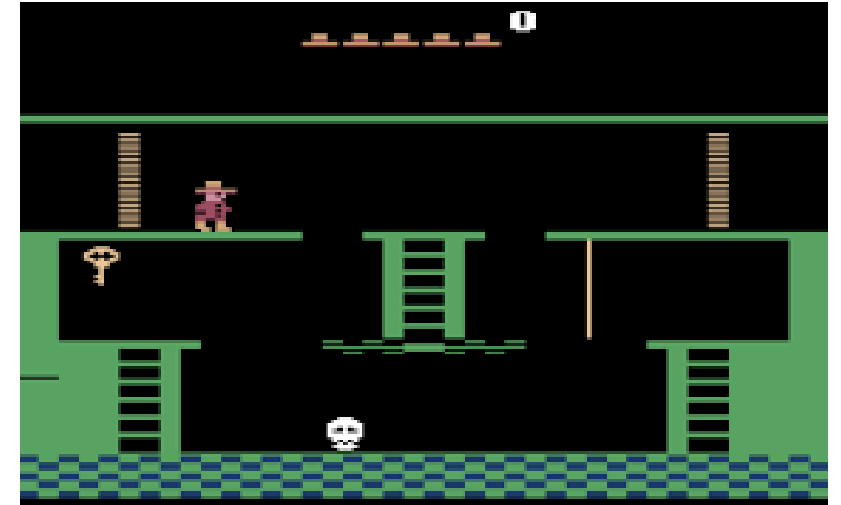
Montezuma's Revenge!

# Exploration

## Bandits

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

- Pseudocounts

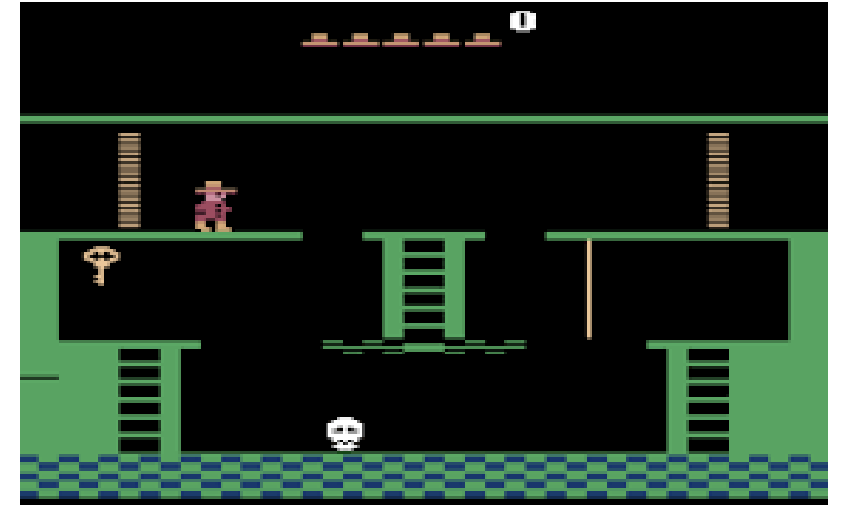


Montezuma's Revenge!

# Exploration

## Bandits

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



Montezuma's Revenge!

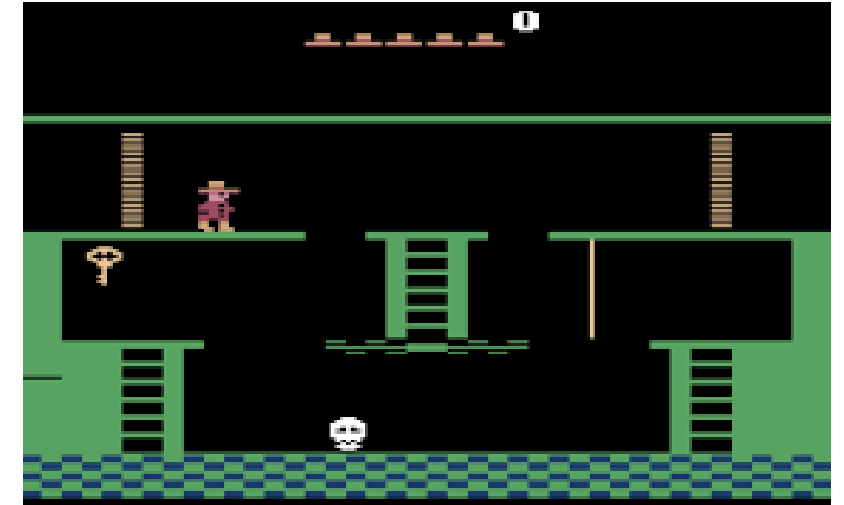
- Pseudocounts
- Curiosity: extra reward for bad prediction

# Exploration

## Bandits

- $\epsilon$ -greedy
- softmax
- UCB  $Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}}$
- Thompson Sampling
- Optimal DP Solution (solving a POMDP!)

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



Montezuma's Revenge!

$$s' = f_{\theta}(s, a)$$
$$f_{\phi}(s, a) = f_{\theta}(s, a)$$



# RL Algorithms

# RL Algorithms

Model  
Based

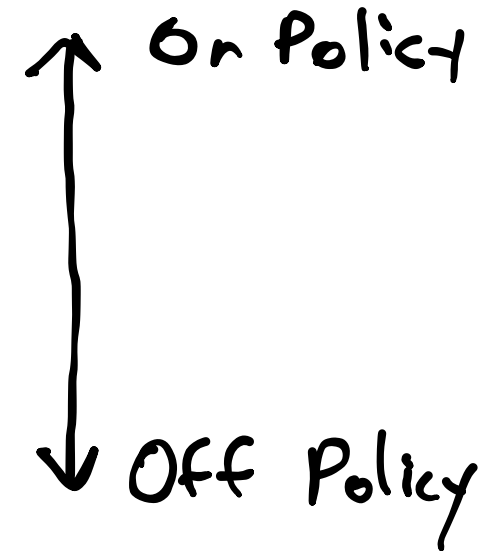
Model  
Free



# RL Algorithms

Model  
Based

Model  
Free



# RL Algorithms

Model  
Based

Model  
Free

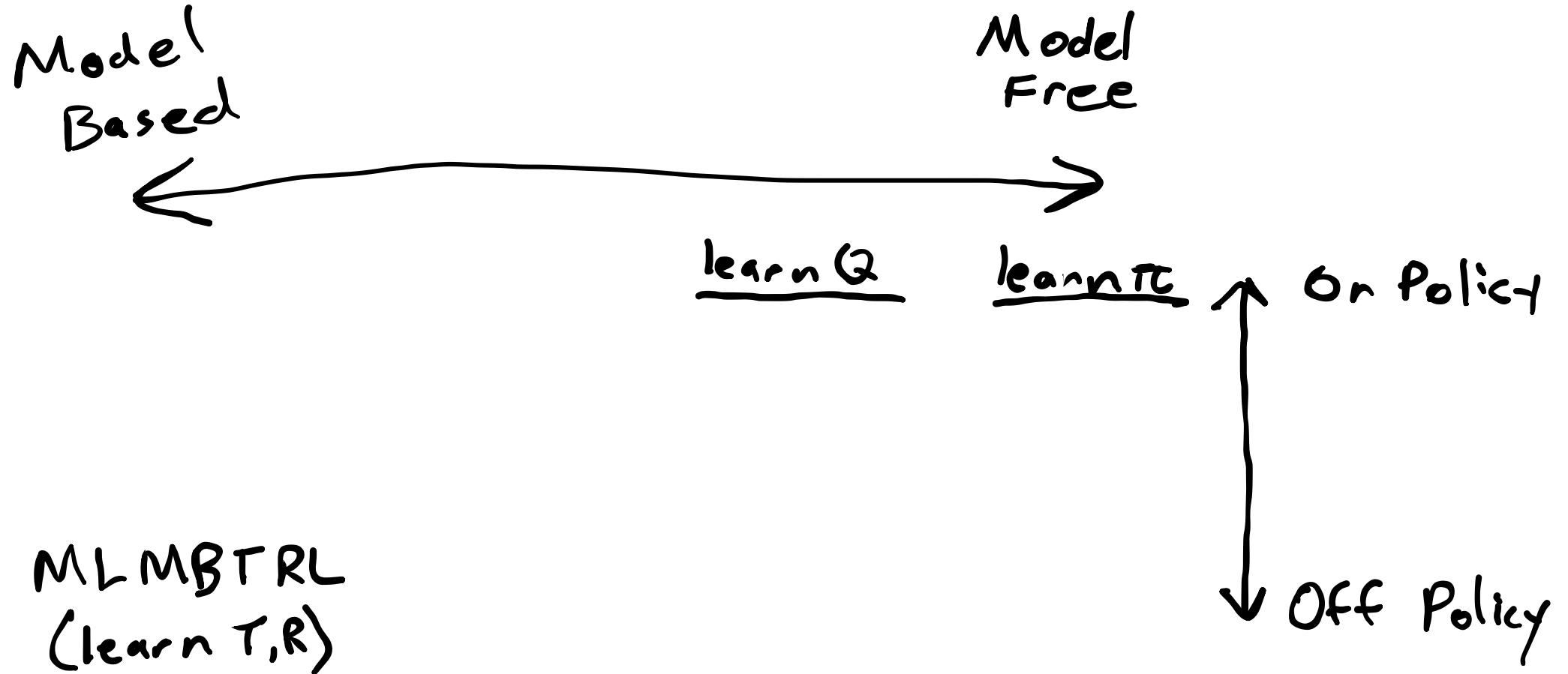


MLMBTRL  
(learn  $T, R$ )

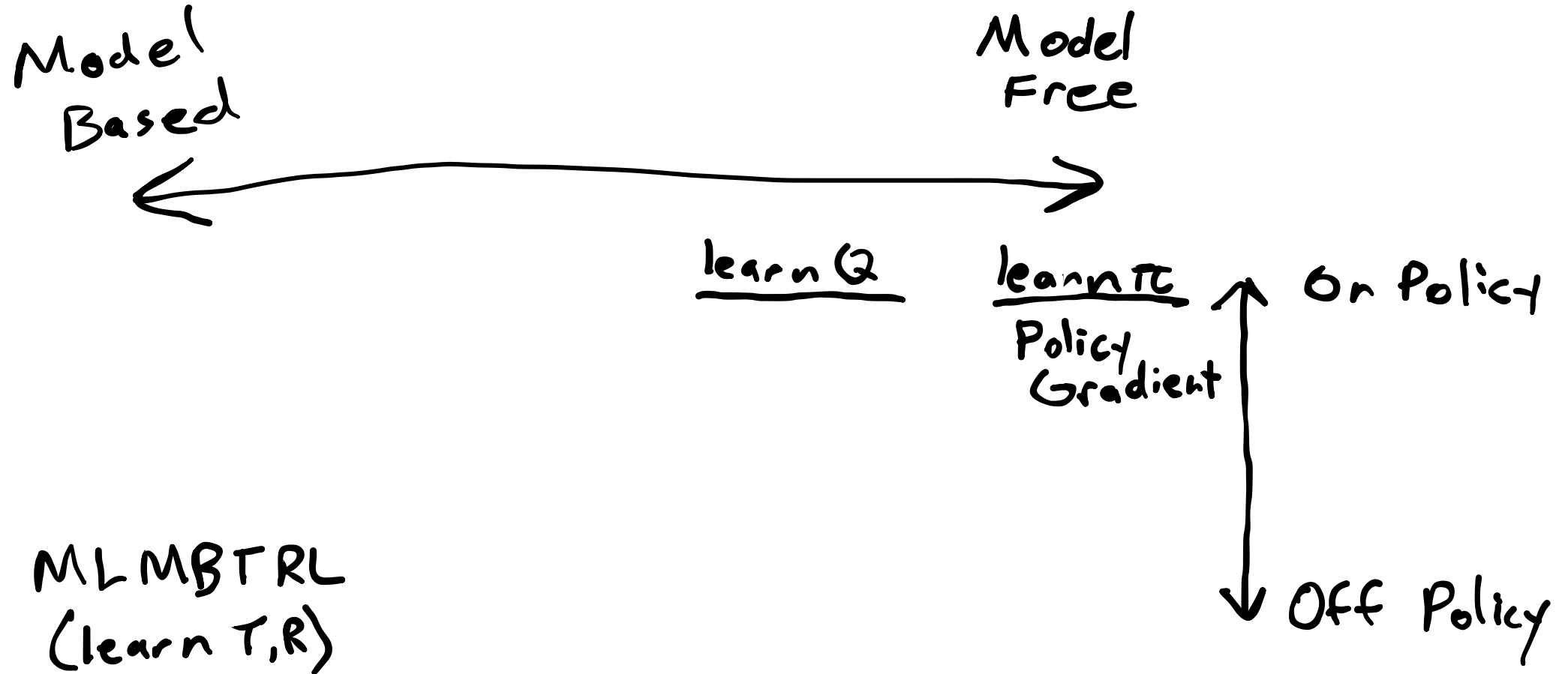
On Policy  
Off Policy

A vertical double-headed arrow pointing from 'On Policy' at the top to 'Off Policy' at the bottom, representing a spectrum of reinforcement learning algorithms based on whether they learn on or off the policy.

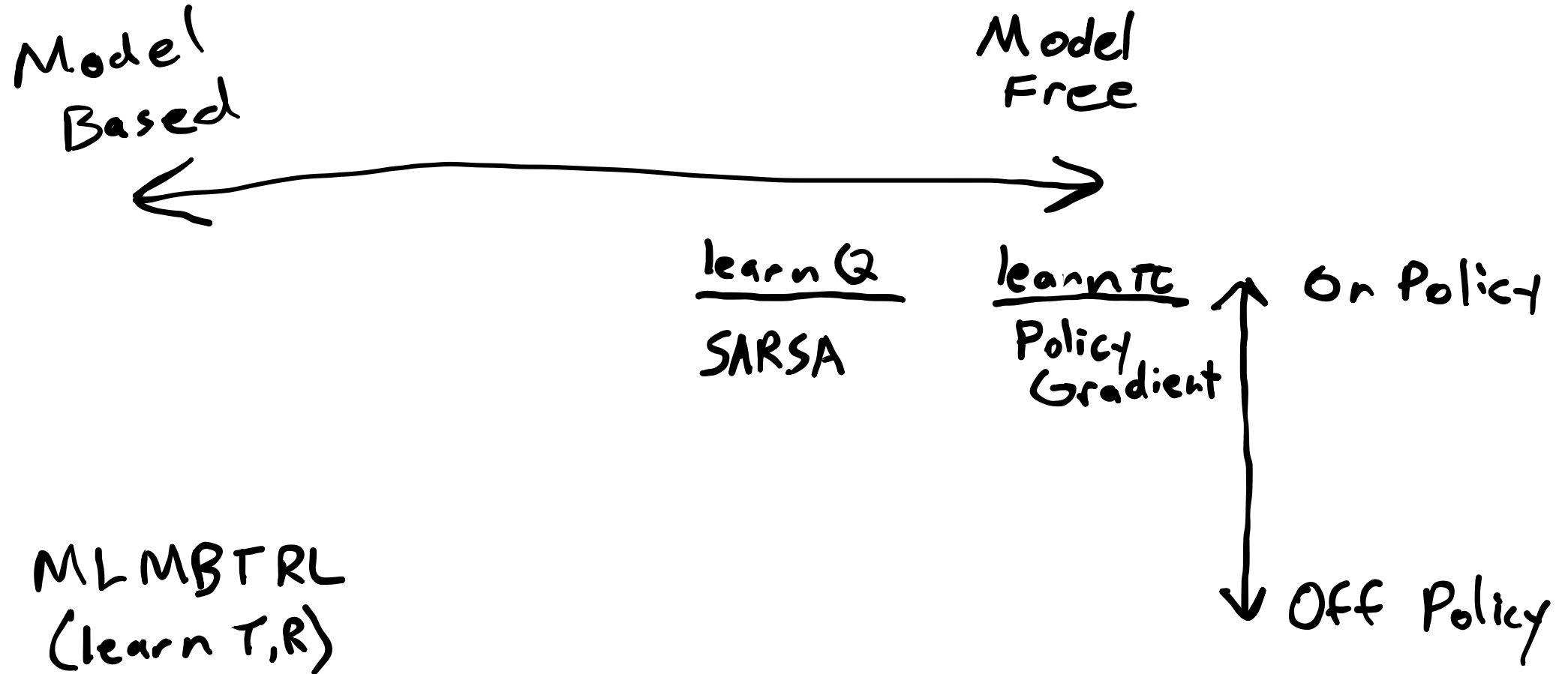
# RL Algorithms



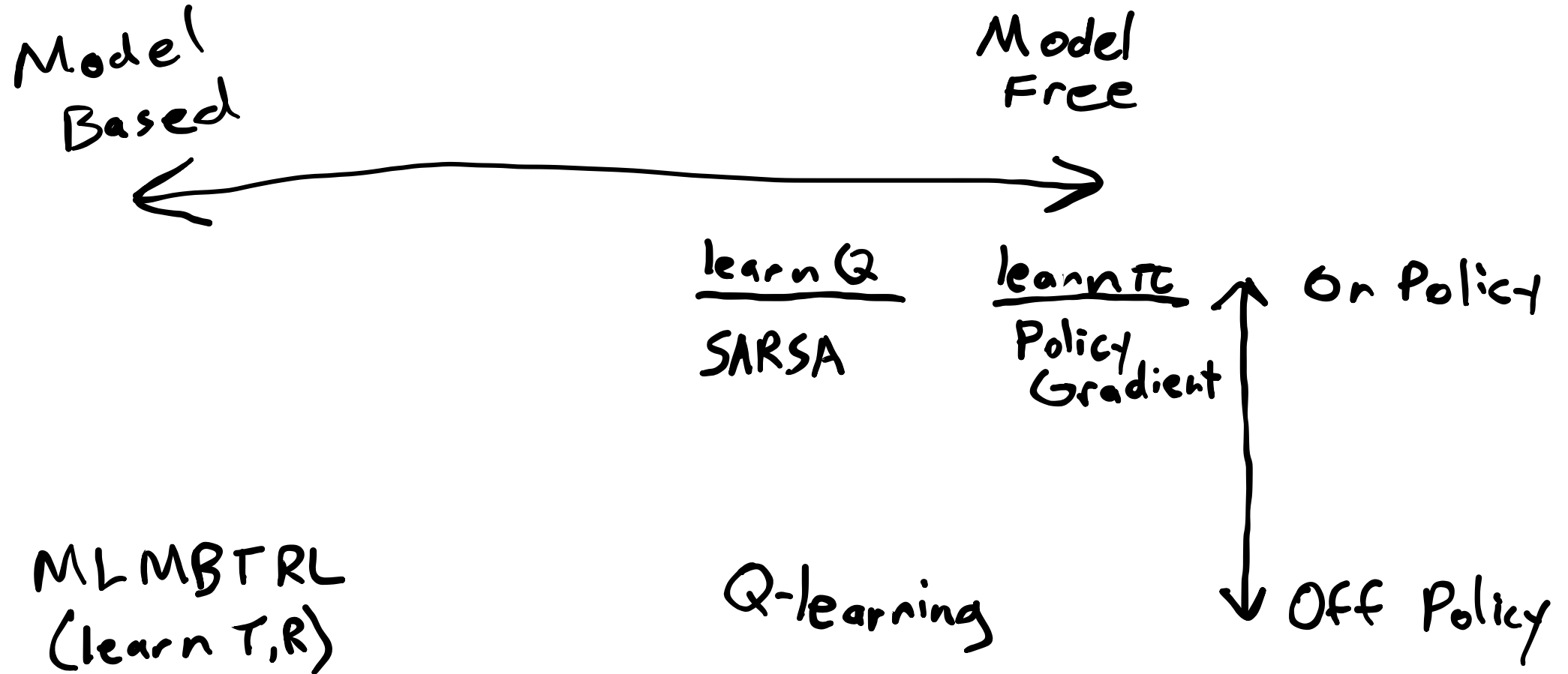
# RL Algorithms



# RL Algorithms

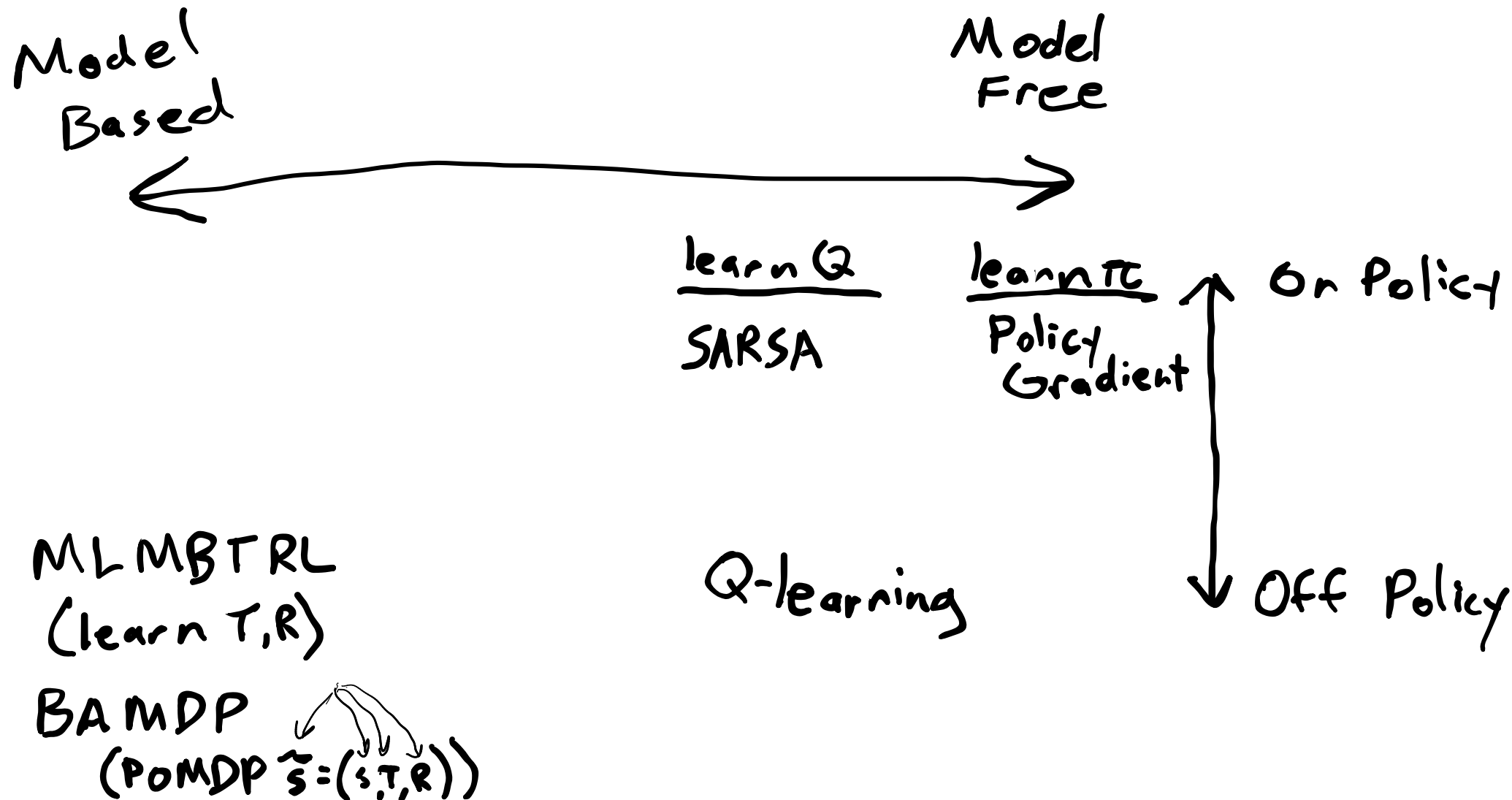


# RL Algorithms





# RL Algorithms



# 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

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

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

↑

# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

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

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

- Natural Gradient

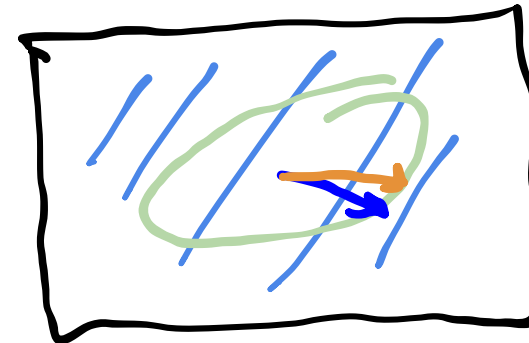
# Policy Gradient

- Likelihood ratio trick
- Causality
- Baseline Subtraction

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

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

- 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

chosen by current policy

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

Eligibility Traces

## Q-learning

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

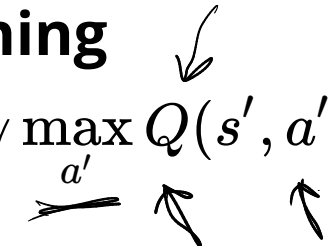
# Q-Learning

## SARSA

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

Eligibility Traces

## Q-learning

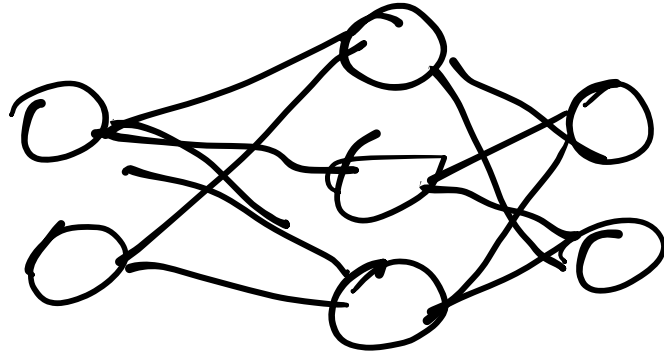
$$Q(s, a) \leftarrow Q(s, a) + \alpha(r_t + \gamma \max_{a'} Q(s', a') - Q(s, a))$$


Double Q Learning

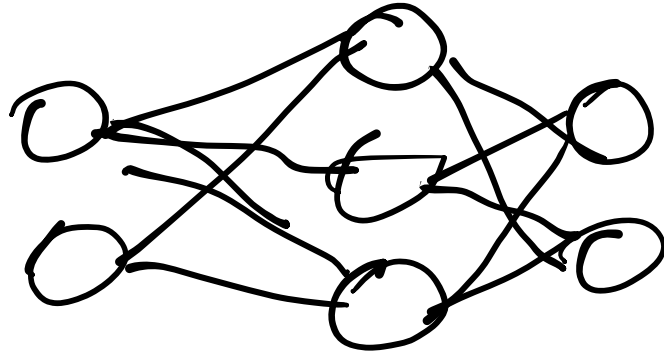


# Neural Networks and DQN

# Neural Networks and DQN



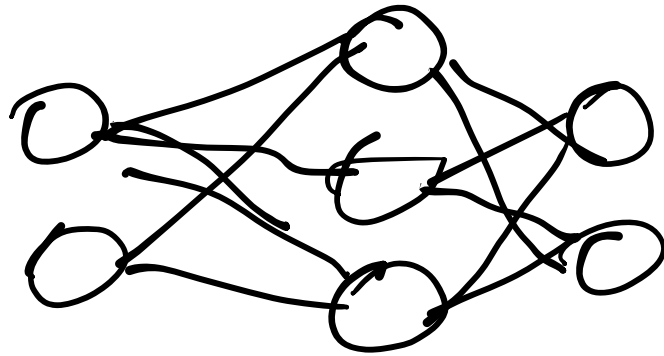
# Neural Networks and DQN



$$f_{\theta}(x) = \sigma(\underbrace{W_2}_{\text{weights}} \underbrace{\sigma(W_1 x + \underbrace{b_1}_{\text{bias}})}_{\text{activation}} + b_2)$$



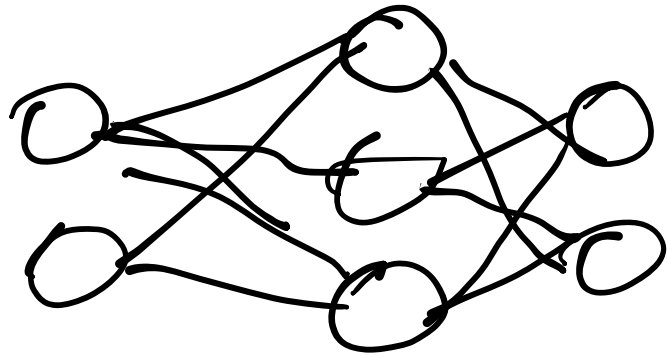
# Neural Networks and DQN



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

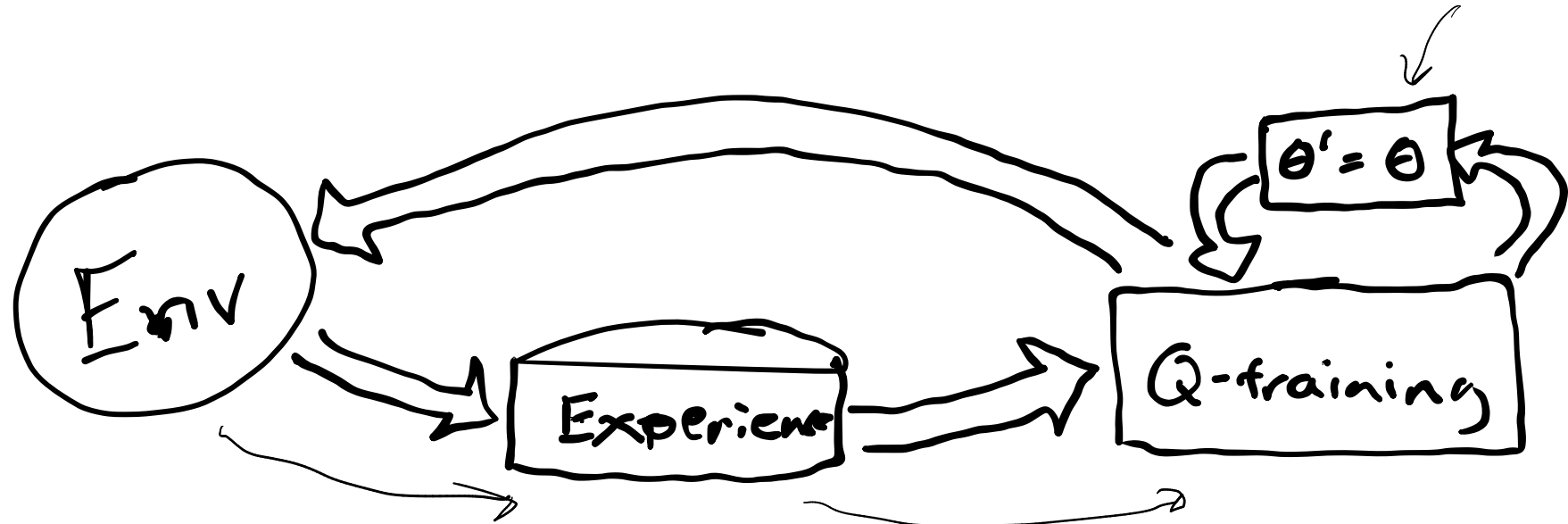
Backprop

# Neural Networks and DQN



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

Backprop



# Actor-Critic

# Actor-Critic

- Actor:  $\pi_{\theta}$

# Actor-Critic

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

# Actor-Critic

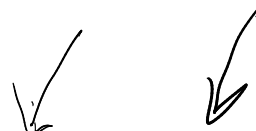
- Actor:  $\pi_{\theta}$
- Critic:  $Q_{\phi}$

## Soft Actor Critic

# Actor-Critic

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

## Soft Actor Critic

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


# POMDPs

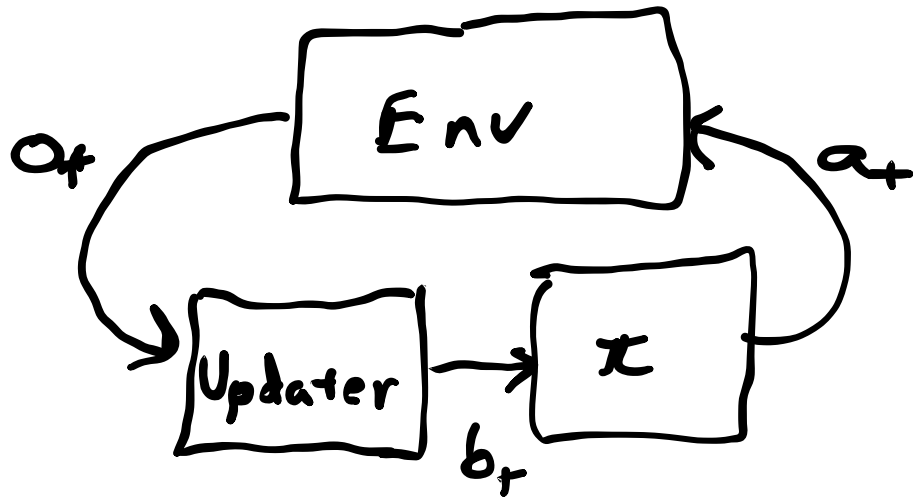


# POMDPs

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

# POMDPs

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

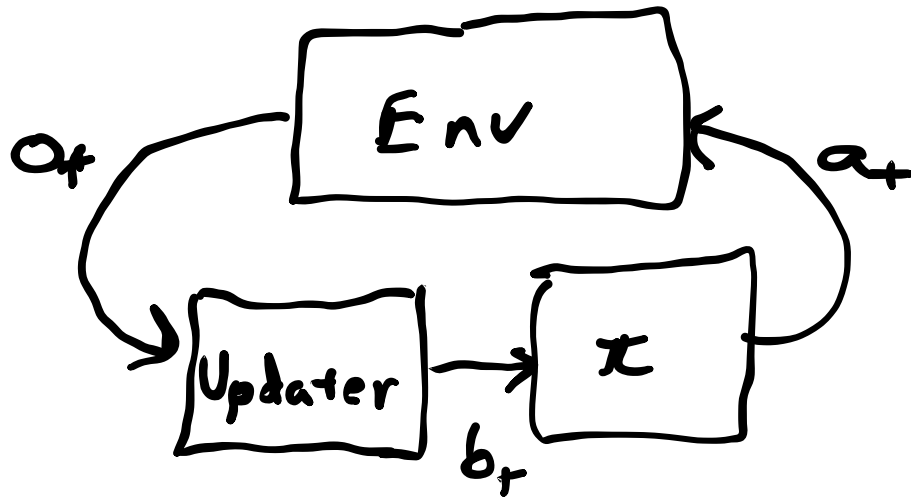


# POMDPs

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

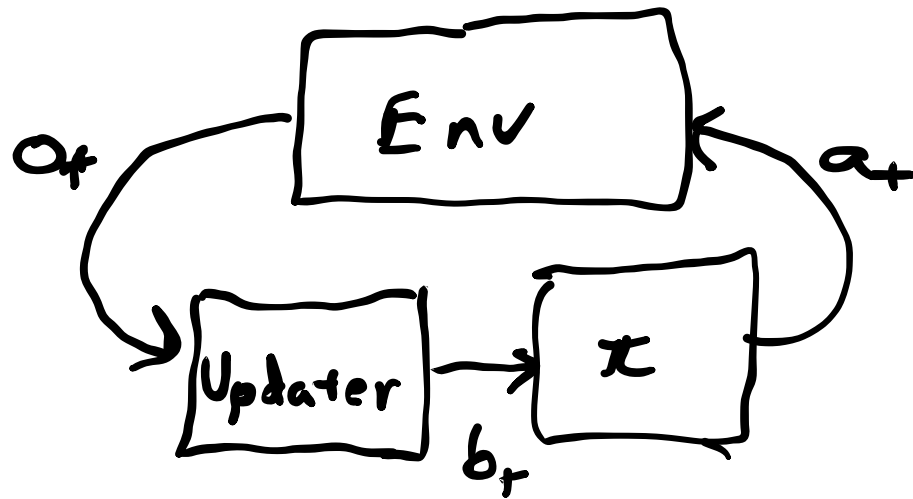
## Belief Updates

- Discrete Bayesian Filter
- Particle Filter



# POMDPs

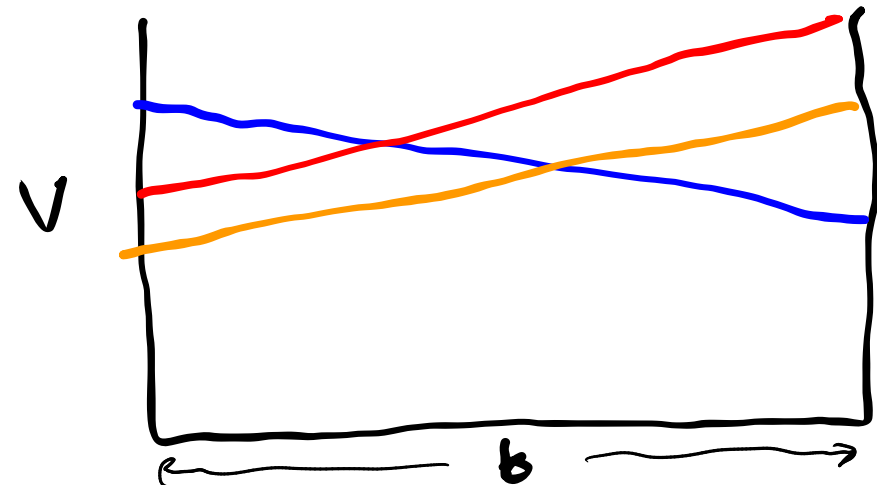
$(S, A, T, R, O, Z, \gamma)$



## Belief Updates

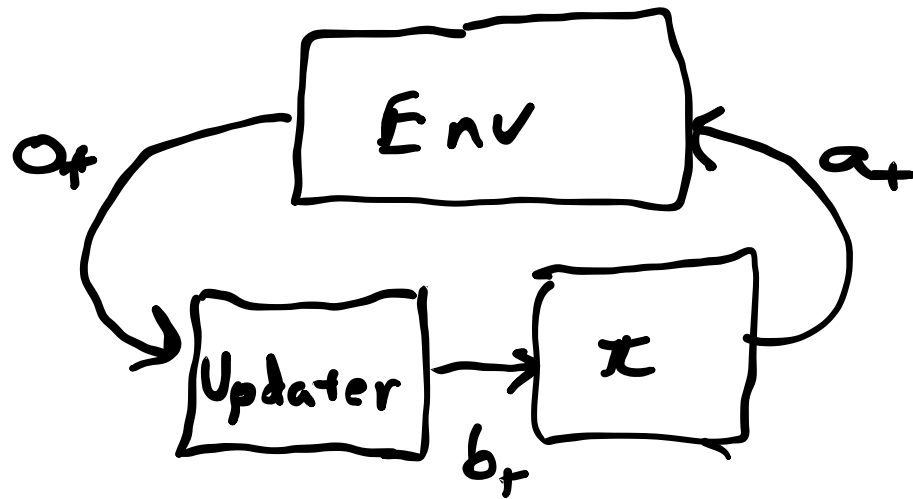
- Discrete Bayesian Filter
- Particle Filter

## Alpha Vectors



# POMDPs

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

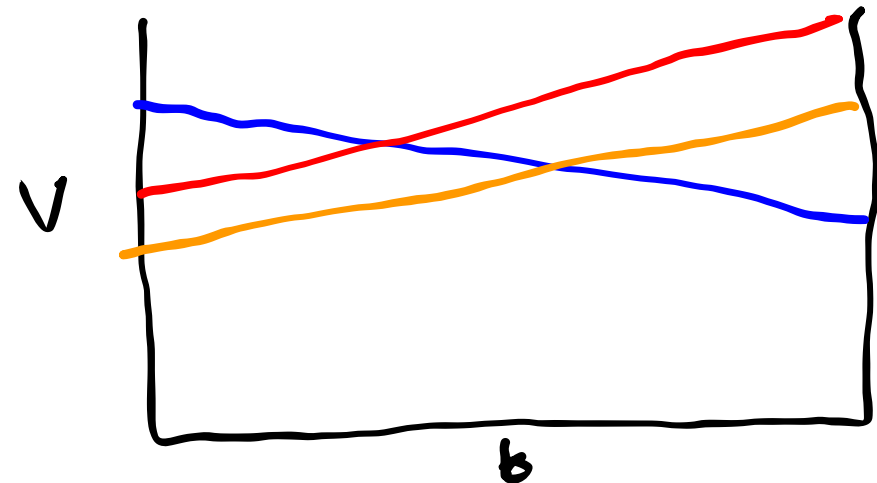


- Each alpha vector corresponds to a conditional plan

## Belief Updates

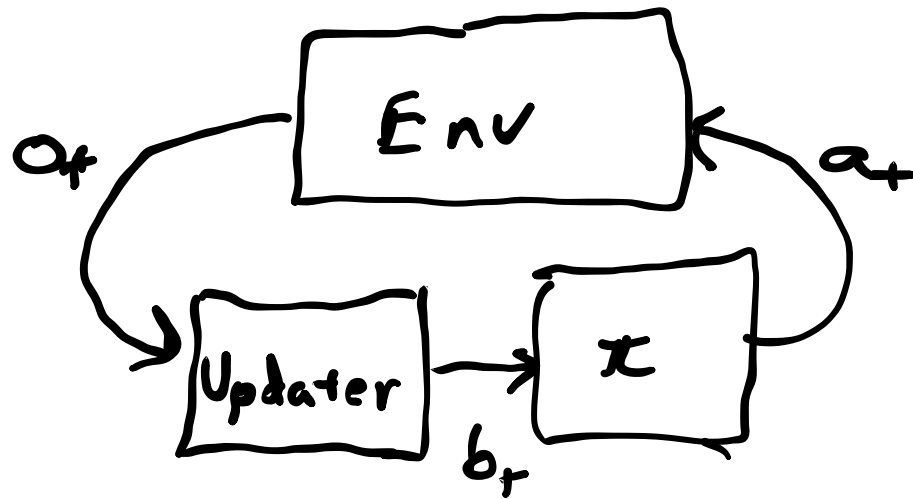
- Discrete Bayesian Filter
- Particle Filter

## Alpha Vectors



# POMDPs

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

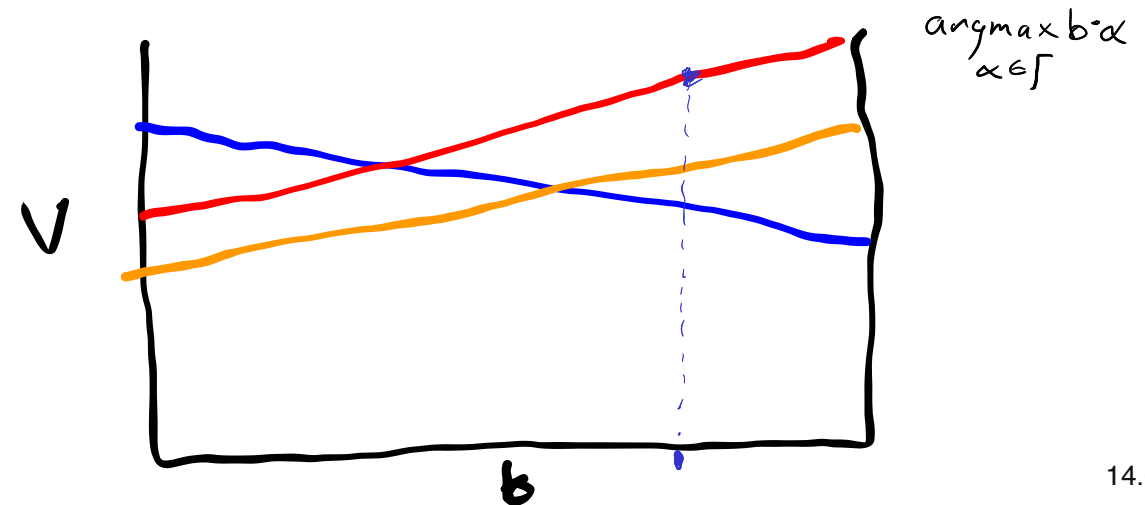


- Each alpha vector corresponds to a conditional plan
- You can prune alpha vectors by solving an LP

## Belief Updates

- Discrete Bayesian Filter
- Particle Filter

## Alpha Vectors



# POMDP Approximations

# POMDP Approximations

## Formulation

- Certainty Equivalence
- QMDP

$$\pi_{MDP}(\underset{s \sim b}{E}[s])$$

$\swarrow$   $\arg\max_s b(s)$

$\nwarrow$  optimal for LQG

$$\pi_{QMDP}(b) = \arg\max_a E_{s \sim b} [Q_{MDP}(s, a)]$$



# POMDP Approximations

## Formulation

- Certainty Equivalence
- QMDP

## Numerical

# POMDP Approximations

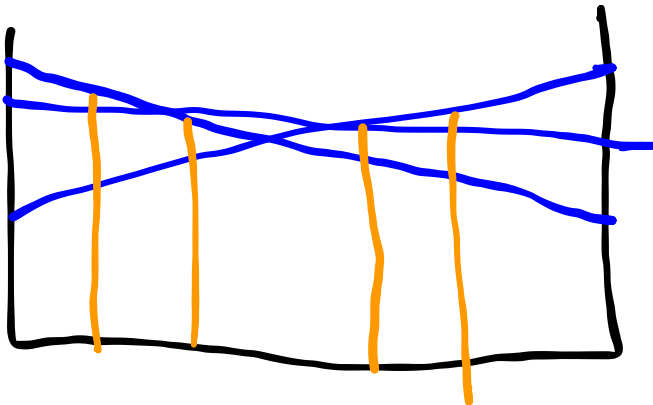
## Formulation

- Certainty Equivalence
- QMDP

## Numerical

### Offline

- Point-Based Value Iteration
- SARSOP



# POMDP Approximations

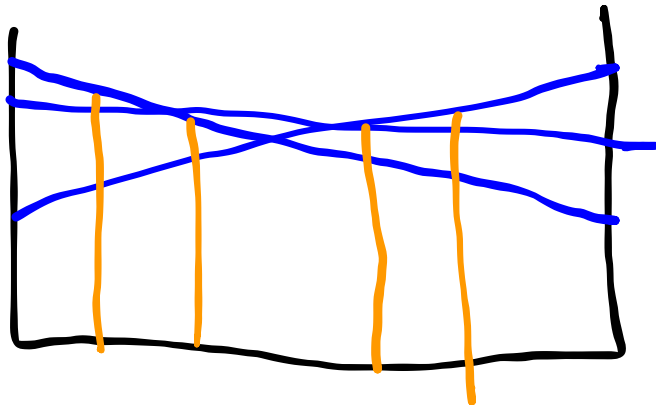
## Formulation

- Certainty Equivalence
- QMDP

## Numerical

### Offline

- Point-Based Value Iteration
- SARSOP



### Online

- POMCP
- DESPOT

# POMDP Approximations

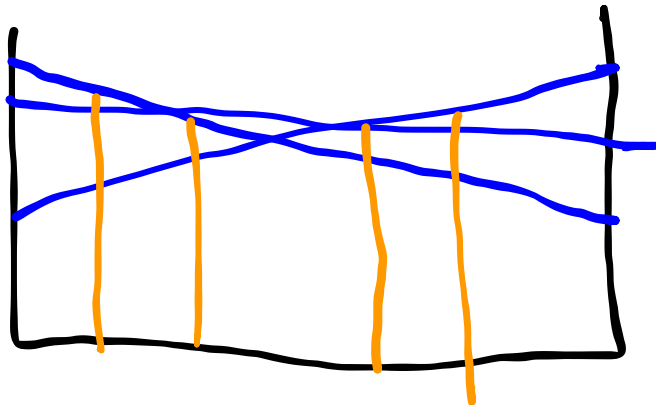
## Formulation

- Certainty Equivalence
- QMDP

## Numerical

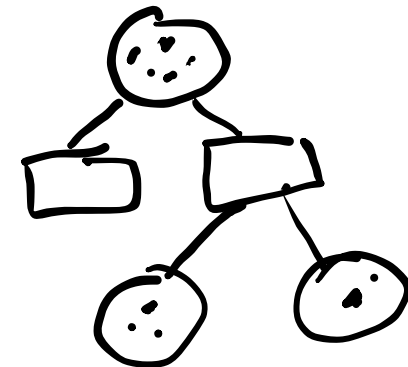
### Offline

- 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)

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

# 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)

up

down

|      |        |        |
|------|--------|--------|
|      | L      | R      |
| up   | -1, -1 | 3, 0   |
| down | 0, -3  | -2, -2 |

NE: every player plays a best response

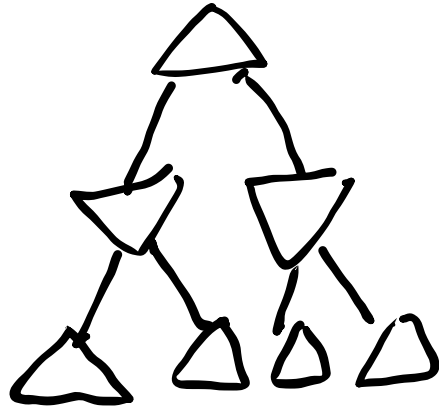
- Every finite game has at least 1 Nash Equilibrium
- Might be pure or mixed
- Algorithms like fictitious play converge in special cases

Dominant Strategy Equilibrium

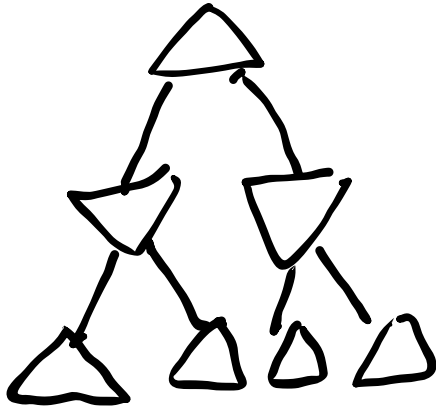
Correlated Equilibrium  $\rightarrow$  NE is a correlated Eq. where all policies are independent

# Turn Taking Games

# Turn Taking Games

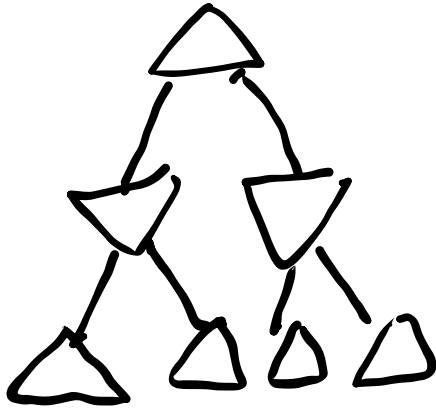


# Turn Taking Games



- Value Function Backup

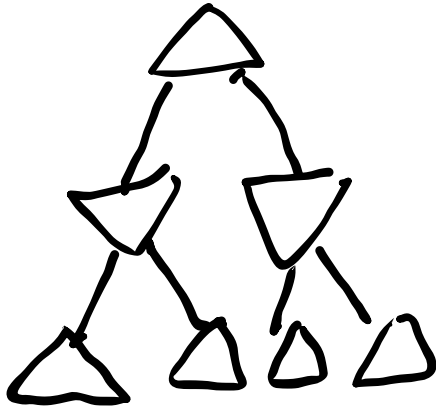
# Turn Taking Games



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

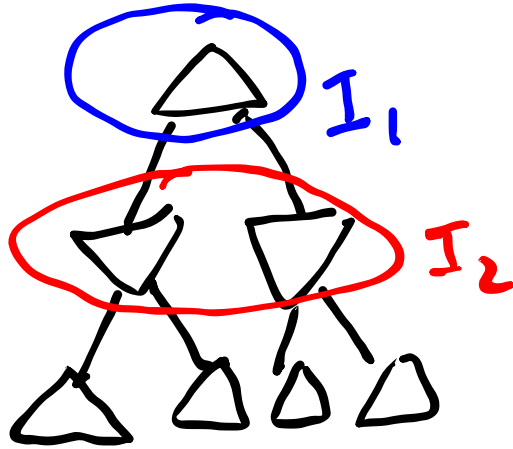


# Turn Taking Games



- 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

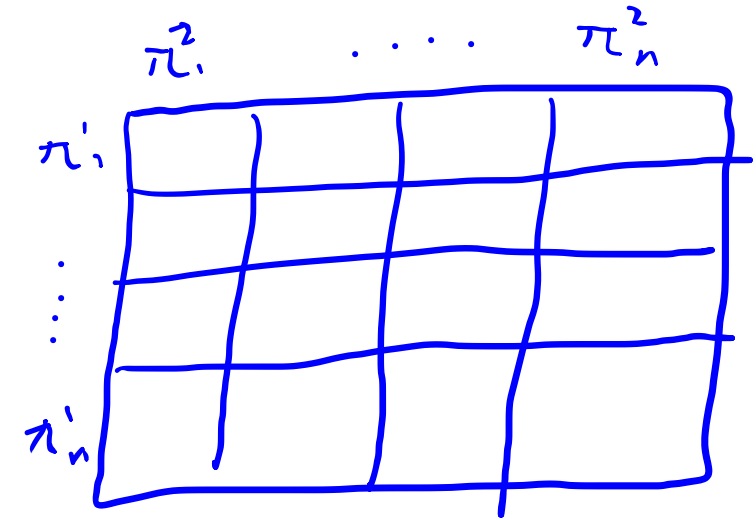
## Markov Games

- 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

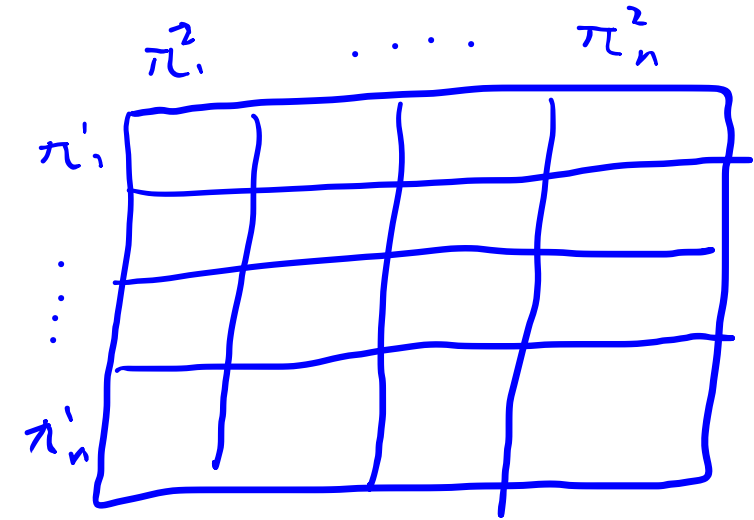
- All players play simultaneously
- Transitions are stochastic
- Best response involves solving an MDP
- Can be reduced to a simple game with policies as actions



# 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



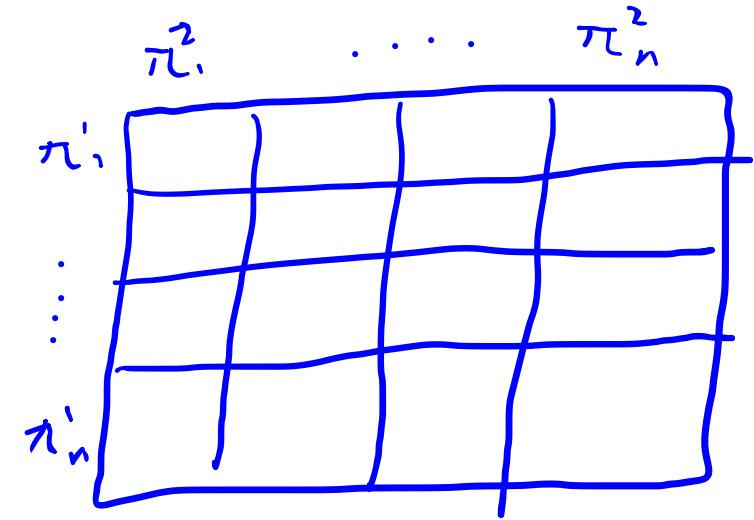
## Partially Observable Markov Games

- 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

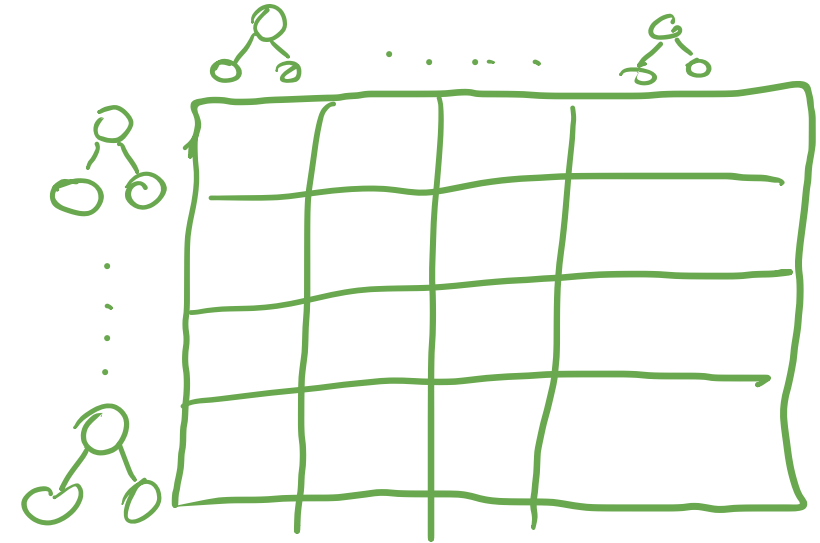
## Markov Games

- 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

initialize  $N(j, a^j, s)$  to 0  
initialize  $\pi^{BR}$

Loop:

1. Simulate MG with  $\pi^{BR}$

2. Update  $N(j, a^j, s) \leftarrow N$  should reflect all past simulations

3.  $\pi^j(a^j | s) \propto N(j, a^j, s) \quad \forall j$

4.  $\pi^{BR} \leftarrow$  best response to  $\pi$  converges to NE for some classes of games

will  $\pi^{BR}$   
converge?  
No

solve an MDP with other player's strategies integrated into it



# 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