

Announcement

ASEN 6519 Advanced Topics
in Sequential Decision Making
Projects due 5pm.

DMU

- Prob. Models
- MDPs
- RL
- POMDPs
- Games

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

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

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

$$2) P(X) = \sum_Y P(X, Y)$$

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

Bayes Rule ✓

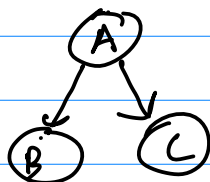
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Independence

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

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

Bayes Nets



$$P(X | \text{all}) = P(X | \text{Pa}(X))$$

$X \perp Y | G$ if all paths between X and Y are d-separated by G

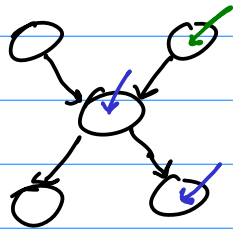
Sampling

Topological Sort then sample from each node

Inference

Given: BN, values of some variables

Output: Distributions of target Variables



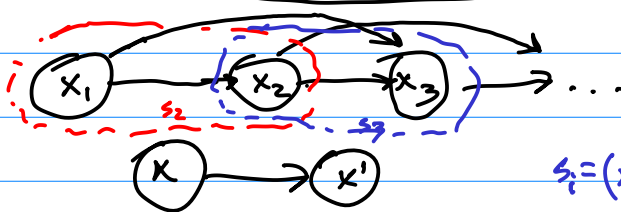
Exact: NP-hard

Approximate :- Direct - Sampling

- Likelihood-weighted Sampling ← Unweighted PF

- Gibbs's method ← Weighted PF

Stochastic Process

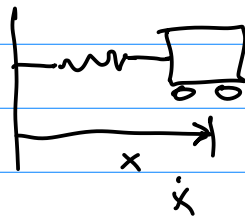


Dyn.
Bayes
Net

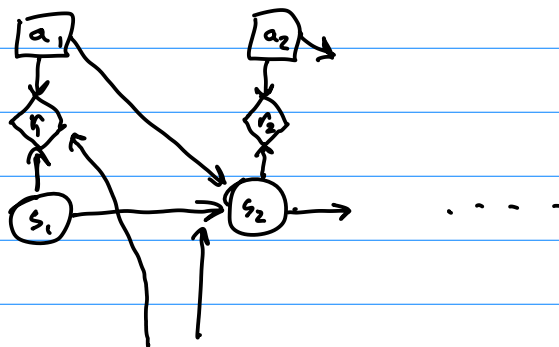
$$s_i = (x_i, x_{i-1})$$

Markov Property

$$P(x_t | x_{t-1}, x_{t-2}, \dots, x_1) = P(x_t | x_{t-1})$$



Decision Network



$s_3 \perp a_1$? No

$s_3 \perp a_1 | s_2$? Yes

$r_k \perp a_1 | s_2 \quad \forall k \geq 2$

$V(s)$

MDP

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

$$S = \{1, 2, 3\}$$

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

$$R(s, a) = E_{s'} [R(s, a, s')]$$

$$T(s' | s, a)$$

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

Bellman's Equation

$$V^*(s) = \max_a \left(R(s, a) + \gamma E[V^*(s') | s, a] \right)$$

\uparrow \uparrow \uparrow

$\sum_{s' \in S} T(s' | s, a) V^*(s')$

Value Iteration

loop

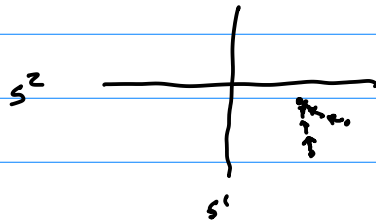
for $s \in S$

$$V^{k+1}(s) = \max_a \left(R(s, a) + \gamma E[V^*(s') | s, a] \right)$$

Converges because Bellman Operator

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

is a contraction in infinity norm



Policy Iteration

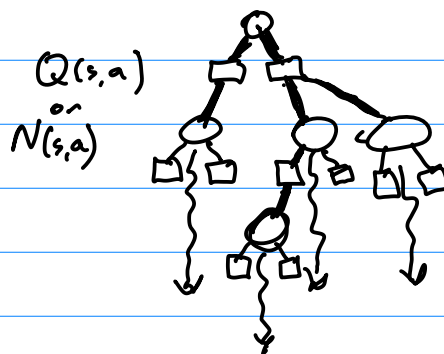
loop

Evaluate Policy

Improve Policy

Online Planning MCTS

Search Expand Rollout Backup



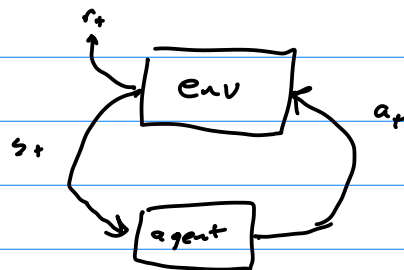
$$Q(s,a) + c \sqrt{\frac{\log N(s)}{N(s,a)}}$$

simulate(...)

$$q \leftarrow r_t + \gamma \text{simulate}(\dots)$$

$$Q(s,a) += \frac{(q - Q(s,a))}{N(s,a)}$$

Reinforcement Learning



$$\theta_1, \theta_2, \dots, \theta_n$$

Exploration vs. Exploitation: Multi-Armed Bandit

ϵ -greedy
softmax
UCB

Thompson Sampling
Optimal DP Solution
 \uparrow POMDP!

$$R(s,a) = \theta_a$$

$$s = (\theta_1, \dots, \theta_n)$$

$$\mathcal{O} = \{0,1\}$$

$$\mathbb{P}(1|a,s') = \theta_a$$

Montezuma's Revenge

Advanced Exploration

- Pseudo Counts $N(s,e)$
- Curiosity : bonus reward if dynamics hard to predict
- Random Network Distillation

RL Algorithms

(s, A, R, γ)

Model Based

MLMBRL
BAMDP

Model-Free

Learn Q

SARSA
On-Policy

Q-Learning
Off-Policy

Learn π

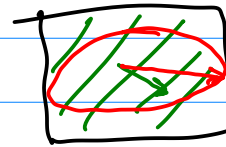
Policy Gradient +

Actor Critic
 π Q, V, A

Policy Gradient

- Likelihood Ratio Trick
- Causality
- Baseline Subtraction
- Natural Gradient

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

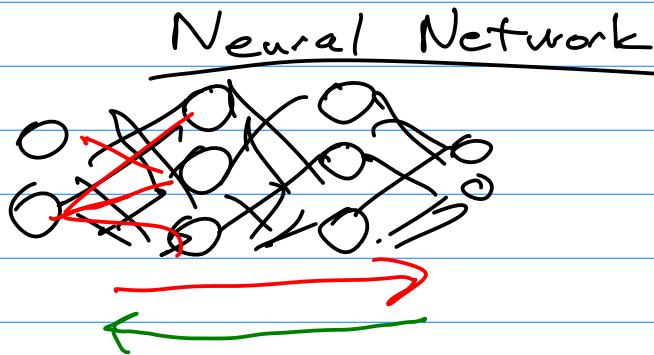


KL-divergence

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

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

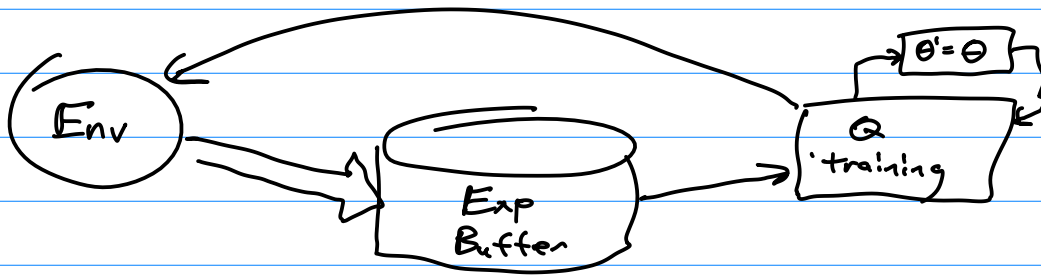
Sarsa ~~≡~~
 Eligibility Traces



$$f_{\theta}(x) = W_3 \sigma(W_2 \sigma(W_1 x + b_1) + b_2) + b_3$$

Back prop

DQN



$Q(s,a) \rightarrow \mathbb{R}$

1. Experience Buffer
2. Periodically freeze target
3. Q-network outputs values for all actions

Rainbow

POMDPs

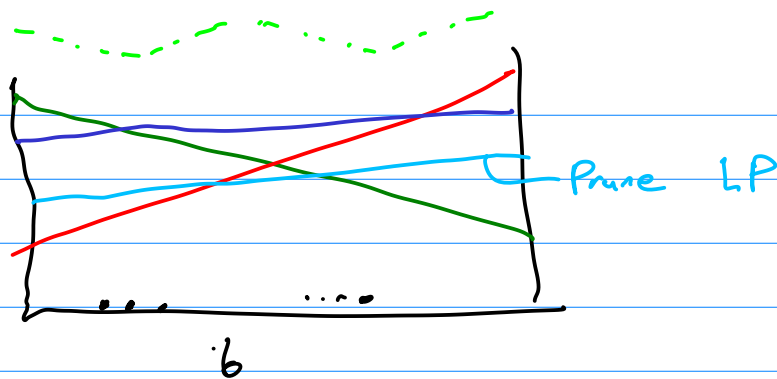
$(S, A, T, R, O, Z, \gamma)$

Belief-Space MDP

Belief-Updates

- Discrete Bayes filter
- Particle Filter

Particle Depletion \wedge



PBVI
SARSOP ←

POMDP is PSPACE
Complete

Formulation Approximation

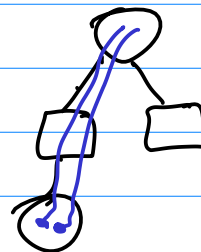
- Certainty Equivalent ← Optimal LQG
- QMDP

Online POMDP Methods

POMCP/PO-UCT/MCTS

DESPOT

↑ determinized scenarios



Games

X Optimal Solution

✓ Nash Equilibria Pure, Mixed

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

Markov Games

POMG / POSG / Extensive Form Games

Imitation

IRL

Transfer Learning

Meta Learning (POMDP)

Successor Features

MAML

θ so that 1 step
get to optimal

Big Problems

	<u>Unc</u>	<u>Tool</u>
1. Immediate + Future Rewards	Outcomp/Action Allegoric	Value
2. Unknown Models	Model/Static	Epistemic RL
3. Partial Observability	State/Dyn.	Epistemic POMDP
4. Other Agents	Interaction Unc.	Games