

Objective of MDP: maximize the expected sum of discounted rewards

Bayes-Adaptive MDP: augment the regular state of the agent with the information it has acquired about the dynamics.

Bayes-optimal behavior:

chooses actions that maximize expected reward as a function of belief-state

learning, planning \Rightarrow planning

environment = MDP = $\langle S, A, T, R, \gamma \rangle$
 $\downarrow \quad \downarrow$
 $S \times A \quad S \times A$

Policy $\pi: S \rightarrow A$

Value V : function of A

optimal policy:

$$\pi^* = \operatorname{argmax}_a Q(s, a)$$

$$Q(s, a) = E[R(s, a)] + \gamma E[V(s')] \quad \left. \begin{array}{l} \\ \\ \end{array} \right\}$$

$$V(s) = \max_a Q(s, a)$$

history H: every finite sequence

of steps in the environment (MDP).

state space of the belief MDP = (belief-state)

1 $S = S \times H$ pairs real states with histories

2 grows exponentially with the length of the amount of history to be considered.

approximate Bayes-optimal behavior:

the value of each action selected is within ϵ of the value of the exact Bayes-optimal action.

Bayesian Sparse Sampling:

preferentially expands only promising parts of the tree by performing rollouts, or simulated trajectories, up to the specified depth.

Sparse Sampling:

recursively expanding a full search tree up to a certain depth d .

- Δ At the root state, each of the A actions is chosen a constant number of times C , yielding a set of $A \cdot C$ children
- Δ Once the tree is fully created, the leaves are each assigned a value of zero.
- Δ Starting at the leaves, the values are backed up and combined via Bellman, giving the parents' values, until the root's value is determined.

Δ Nodes in the search tree = $(A \cdot C)^d$

goal: addressing exploration vs exploitation dilemma

policy: action value function.
(Q function)
 $Q(s, a)$

Q function satisfies Bellman equation.
 $Q(s_t, a_t) = E(r_t | s_t, a_t) + \gamma E[\max_{a \in A} Q(s_{t+1}, a) | s_t, a_t]$

RL: learning policy.

② Value based: (learning + planning)
estimate the optimal Q-function directly Q-learning from which a greedy policy is recovered.

policy based: (learning + planning)
estimate a good policy directly.

model based: (planning)
estimate the transition and reward models, and then determines a policy by solving the planning problem in the learned model

Bayesian RL: model based.

a prior distribution is defined over transition and reward models: $P(\theta, \mu | s_0)$

given experience (history)

$s_0 a_0 r_0 s_1 \dots s_t a_t r_t s_{t+1}$

one determines the posterior distribution

$P(\theta, \mu | s_0 a_0 r_0 s_1 \dots s_t a_t r_t s_{t+1})$

learning: updating the posterior

learning \Rightarrow planning (action selection)

A Bayesian approach to learning optimally in a MDP is essentially equivalent to solving a POMDP.

Classical action selection:

Δ ϵ -greedy: with probability $1-\epsilon$ choose the current best estimate $a^* = \arg \max_a \hat{Q}(s, a)$. Otherwise choose a random action $a \in A$

Δ Boltzmann: Sample a random action according to $P(a|s) = \frac{\exp(\hat{Q}(s, a)/\tau)}{Z}$ where τ is a temperature parameter and Z is a normalization constant.

Δ Interval estimation: Choose an action according to

$a = \arg \max_a [\hat{Q}(s, a) + U(s, a)]$

where $U(s, a)$ is a $(1-\delta)$ upper confidence interval on the point estimate $\hat{Q}(s, a)$.

belief-state MDP

||

Bayes adaptive MDP.

||

meta-level MDP

state: b_t is given by a current

base-level state s_t and

a posterior distribution over the base-level transition model θ and reward model μ .

$b_t = \langle P_t^\theta P_t^\mu s_t \rangle$

$\begin{cases} P_t^\theta = P(\theta | s_0 a_0 \dots s_{t-1} a_{t-1} s_t) \\ P_t^\mu = P(\mu | s_0 a_0 r_0 \dots s_{t-1} a_{t-1} r_{t-1}) \end{cases}$

③ meta-level states

$$b_t = \langle P_t^\theta P_t^\mu S_t \rangle$$

||
histories

$$b_t \equiv S_0 a_0 r_0 \dots S_{t-1} a_{t-1} r_{t-1} S_t$$

state transition probability

the probability of a particular history extension r_t, S_{t+1} given the current history

$S_0 a_0 r_0 \dots S_{t-1} a_{t-1} r_{t-1} S_t$ and action a_t .

meta-level belief states =

base-level histories

action selection strategy in Bayes-RL setting (myopic).

Δ Thompson Sampling:

given a current belief state $b_t = \langle P_t^\theta P_t^\mu S_t \rangle$.

sample a transition and reward model, θ and μ , from the belief state distributions

P_t^θ and P_t^μ , solve for the optimal Q-function

$Q_{\theta\mu}(s, a)$ for this model, then select

the optimal action $a_t = \underset{a \in A}{\text{argmax}} Q_{\theta\mu}(S_t, a)$.

action selection strategy in

Bayes-RL setting (optimal).

enumerating possible futures, averaging according to their realization probabilities, and choosing the best action.

the only guaranteed way to approximate Bayes optimal action selection at a given belief state is to simulate the belief state MDP to the effective horizon.

reward:

$$P(r_t | P_t^\theta P_t^\mu S_t a_t) = \int P(r_t | S_t a_t, \mu) P_t^\mu(\mu) d\mu$$

$$= P(r_t | S_0 \dots S_t a_t)$$

transition:

$$P(P_{t+1}^\theta P_{t+1}^\mu S_{t+1} | P_t^\theta P_t^\mu S_t a_t)$$

$$= \frac{1}{\int P_{t+1}^\theta = P(\theta | S_0 a_0 \dots S_t)}$$

$$\int P_{t+1}^\mu = P(\mu | S_0 a_0 \dots S_t) \left[\int P(S_{t+1} | S_t a_t, \theta) P_t^\theta(\theta) d\theta \right]$$

$$\int_{r_t} \int_{\mu} \frac{1}{\int P_{t+1}^\mu = P(\mu | S_0 \dots S_t a_t r_t)} P(r_t | S_t a_t, \mu) P_t^\mu(\mu) d\mu dr_t$$

$$= P(r_t S_{t+1} | S_0 \dots S_t a_t)$$

Paper:
Support planning in domains
with richly structured prior knowledge.

$$Q^*(\langle s_t, h_t \rangle, a) = \max_{\pi} E_{\pi} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'} \mid a_t = a \right]$$

MDP: 5-tuple

$$M = \langle S, A, P, R, \gamma \rangle$$

S: states

A: actions

P: $S \times A$: transition probability.

R: $S \times A$ bounded reward function.

γ : discount factor

assume ~~P~~ unknown ("dynamics").

P is a latent variable subject to $P(\phi)$

observing a history of actions and states

$$h_t = s_1 a_1 s_2 a_2 \dots a_{t-1} s_t$$

posterior belief on ϕ updated using

Bayes' rule.

$$P(\phi | h_t) \propto P(h_t | \phi) P(\phi).$$

(augmented transition)

$$P^+(\langle s, h \rangle, a, \langle s', h' \rangle) = \mathbb{1}[h' = h, a, s'].$$

(transition)

(dynamics)

$$\int P(s, a, s') P(\phi | h) d\phi$$

$$R^+(\langle s, h \rangle, a) = R(s, a).$$

(augmented reward)

augmented MDP: Bayes-adaptive MDP =
belief state MDP: meta-level MDP.

$$M^+ = \langle S^+, A, P^+, R^+, \gamma \rangle$$

dynamics of the BAMDP are known.

(P^+).

1. Use UCT to allocate search effort to promising branches of the state-action tree
2. use sample-based rollouts to provide value estimates at each node.
3. only sample a single transition model P^i from the posterior at the root of the search tree at the start of each simulation i .
4. use P^i to generate all the necessary samples during this simulation.

essentially BAMDP is a POMDP

~~this root~~
★ state node:

$$\text{belief states: } \langle s, h \rangle$$

★ action nodes: a .

★ visit counts: $N(\langle s, h \rangle)$ for state nodes
 $N(\langle s, h \rangle, a)$ for action nodes

initialized to 0.
updated throughout search.

★ Value. $Q(\langle s, h \rangle, a)$ initialized to 0.
maintained for each action node.

★ each simulation traverses the tree without backtracking by following the UCT policy.

at state nodes defined by

$$\arg \max_a Q(\langle s, h \rangle, a) + c \sqrt{\frac{\log(N(\langle s, h \rangle))}{N(\langle s, h \rangle, a)}}$$

↓

$a \Rightarrow$ traverse the tree

★ given an action, the transition distribution p_i corresponding to the current simulation i is used to sample the next state.

if the transition parameters for (θ, ϕ) different states s and actions a are independent, we can completely forgo sampling a complete \mathcal{P} , and instead draw any necessary parameters individually for each state-action pair.

measure metrics in evaluation:

total undiscounted reward over many steps
while the paper optimizes for the discounted reward from the start state

Lazy Sampling

Sample \mathcal{P} lazily, creating only the particular transition probabilities that are required as the simulation traverses the tree, and also during the rollout.

$\mathcal{P}(s, a, \cdot)$ to be parametrized by a

latent variable $\theta_{s,a}$ for each state and action pair.

depending on an additional set of latent variables ϕ .

$$P(\theta | h) = \int_{\phi} P(\theta | \phi, h) P(\phi | h) d\phi.$$

$$\theta = \{\theta_{s,a} | s \in S, a \in A\}$$

$$\theta_t = \{\theta_{s_1, a_1}, \dots, \theta_{s_t, a_t}\}$$

$$P(\theta | \phi, h) = P(\theta_{s_1, a_1} | \phi, h) P(\theta_{s_2, a_2} | \theta_1, \phi, h)$$

$$\dots P(\theta_{s_T, a_T} | \theta_{T-1}, \phi, h) P(\theta | \theta_T | \theta_T, \phi, h).$$

results

△ planning time - performance trade-off

△ performance scale as a function of planning time.

Compatible with arbitrary priors (\mathcal{P}).

Motivation:

Solve large tasks with uncertainty about the dynamics

want to take advantage of structured prior knowledge (sample latent param)

rewards with discounts.

"transitions correlated"



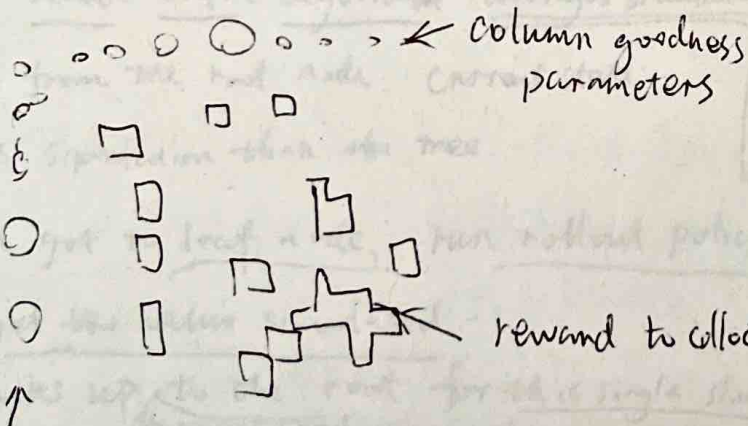
can play in

arbitrary prior
to it.

formalism: Bayesian model-based RL
distribution over R, P

evaluate: infinite grid task.

2-D



row goodness parameters.

Typical MDP = $M = \langle S, A, P, R, \gamma \rangle$

but P is a latent variable with prior $P(P)$ (arbitrary, / prior knowledge)

Goal: Find exploration policy that maximizes expected sum of discounted rewards

$$\int_{\mathcal{P}} P(P) \mathbb{E}_{\pi(P)} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, h_0 = s \right] dP \quad \text{over all possible dynamics that can happen.}$$

all possible transitions according to your posterior

Equivalent to solving an augmented MDP in belief space with known dynamics (\mathcal{P})

$$p^+(\langle s, h \rangle, a, \langle s', h' \rangle) = \mathbb{1}(h' = has) \int_{\mathcal{P}} p(s, a, s') P(P) dP$$

Previously: { don't scale
don't work with any prior
direct or multinomial models

Major obstacle:
Computationally intractable to solve.

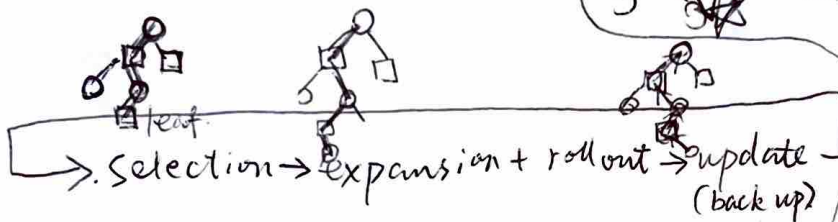
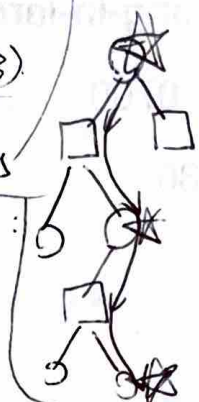
except for bandits (MAB)

{ state space
posterior transition states

info optimal BOP policy

Kocsis & Szepesvári work:

MCTS search



forward search online algorithm. Averages simulation

- Δ Start from the root node: current state
- Δ Run simulation thru the tree
- Δ when get to leaf node, run rollout policy to get the value simulated.

Δ backs up to the root for this single simulation.

Δ average over the nodes on the path. the reward of nodes on path
simulations

UCB: upper confidence bounds.

UCT: UCB applied to trees

In

one simulation

state node

Δ treat each state node as a multi-armed bandit problem (use UCB)

state node decides where to go in the tree.

put more effort to where it seems fruitful.

(can see most of the rewards)

Solution 1:

Bayesian adaptive

BA-UCT: leverage modern

Sample-based MDP Solver:

MCTS/UCT applied to BA-MDP.

Issue: expensive belief updates at every tree node.

multiple ~~every~~ simulation, each with integrating over posterior and sample from ^{transition} posterior. belief update

100x simulations for each

needs to run many simulations
not feasible for generating priors
Monte Carlo.

Solution 2:

BA-UCT + Root Sampling

Silver & Veness 2012

(belief update)

Sample p at root for each sim

one belief update per step

still Bayes-optimal in the limit

Use this single p (sampled dynamics) throughout the sim.

Issue:

wasteful sampling, each sim only requires a small subset of all parameters of p .

[factorization of the prior]

develop some conditional independence

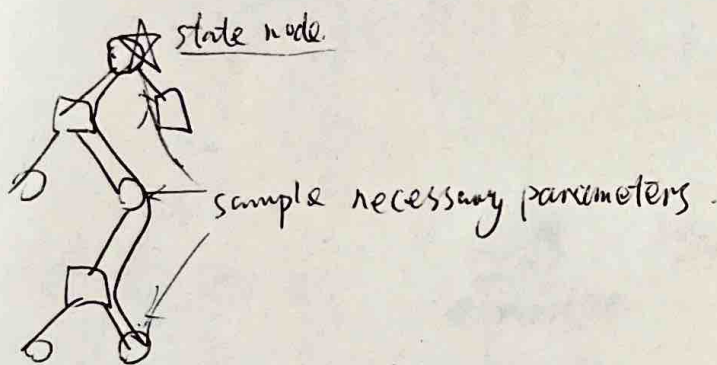
[only sample parameters we need]

Solution 3

BA-UCT + RS + Lazy Sampling (LS):

Sample MDP param on demand
based on current simulation.

Idea: ✓ minimize sampling overhead, put more effort on search.



③ standard domains

- double-loop
- Grid 5
- Grid 10
- Dearden Maze.

BAMCP get best sum of rewards despite the fact that it optimizes for discounted rewards

~~finite~~ sum of infinite.

planning time ↑, return — or ↓

for other algo

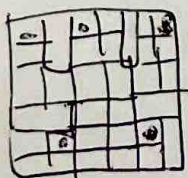
because of over ~~optimistic~~ optimistic

BA-UCT + RS + LS + Rollout Learning
learn better rollout policy online.



Bayes-adaptive Monte Carlo Planning.

eval:



1-step rollout: choose action $\pi_k(s_t)$ and simulate the system to get reward r_t , then we use this reward to update our estimate of the value function: $V_{k+1}(s_t) = (1-\gamma)V_k(s_t) + \gamma \cdot X \cdot r_t$

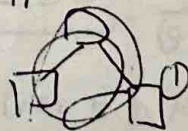
Dearden's Maze.

{ 264 states

{ sparse Dirichlet-Multinomial prior

Appendix:

Solution 2:



① set up the MCMC at the root node.

② feed samples ^{of P} to UCT algorithm.

planning time ↑, return ↑

prior of transition distri matters a lot.

prior affects behavior and performance.
(return, reward)