

recap

Markov Decision Processes useful framework for describing sequential decision making under uncertainty
if MDP is given (know its transition & reward models as well as its state & action space)
and its not too "big" (small, discrete state & action spaces)
then can use dynamic programming / Bellman operator to do value iteration
when the MDP is unknown and/or large, need additional techniques
if we have a good estimate of optimal V^* can use to obtain good control performance

Problem Setting

batch set of data D , consisting on N trajectories of length H

$s_{11}, a_{11}, r_{11}, s_{12}, a_{12}, r_{12} \dots$
time index into traj

gen by behavior π_b
(could be nonstationary,
unknown, multiple)

S is very large or continuous, may be represented as a
vector of feature values

goal: estimate optimal state action values * relates to 1st
main challenge

Why is this an interesting, important setting?

- in many real world applications (customer marketing,
electronic health records, intelligent tutoring systems)
have access to prior data about seq of decisions and
their outcomes

- but may be hard/costly to gather more data

- would like to use existing data to inform future decision
making

- Why large/factored S ? realistic
 \check{V} describe

- Why estimate Q ?

sufficient for making good decisions

more precisely, if compute a close estimate to V^*/Q^*
then greedy policy wrt that estimate will have

close to V^* performance in real world (at least in discrete SA)

thm: (Satinder & Yee 1994). Let V^* be the optimal
value func for a discrete-time MDP having finite
state and action sets and an ∞ horizon w/ geometric
discounting $\gamma \in [0,1]$. If \tilde{V} is a function s.t. $\forall s \in S$

$$|V^*(s) - \tilde{V}(s)| \leq \epsilon$$

and the greedy policy $\pi_{\tilde{V}}$ is defined as

$$\pi_{\tilde{V}}(s) = \arg \max_{a \in A(s)} r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) \tilde{V}(s')$$

then for $\forall s \in S$

$$V^*(s) - V^{\pi_{\tilde{V}}}(s) \leq 2\gamma\epsilon/(1-\gamma)$$

define what
 $\pi_{\tilde{V}}$ means

proof:

for shorthand define loss $L_{\tilde{V}}(s) = V^*(s) - V^{\pi_{\tilde{V}}}(s)$

\exists state z that achieves the max loss. Call this state z

$$\forall s \in S, L_{\tilde{V}}(z) \geq L_{\tilde{V}}(s)$$

for state z consider an optimal action $a = \pi^*(z)$ and

the action specified by $\pi_{\tilde{V}}$, $b = \pi_{\tilde{V}}(z)$.

because $\pi_{\tilde{V}}$ is the greedy π for \tilde{V} , b must look at least as good

$$r(z, a) + \gamma \sum_{s' \in S} p(s'|z, a) \tilde{V}(s') \leq r(z, b) + \gamma \sum_{s' \in S} p(s'|z, b) \tilde{V}(s')$$

proof cont

because $\forall s' \in S \quad V^*(s) - \epsilon \leq \tilde{V}(s') \leq V^*(s') + \epsilon \quad (\text{by assumption})$

$$r(z, a) + \gamma \sum_{s' \in S} p(s'|z, a) (V^*(s') - \epsilon) \leq r(z, b) + \gamma \sum_{s'} p(s'|z, b) (V^*(s') + \epsilon)$$

$$(*) \quad r(z, a) - r(z, b) \leq 2\gamma \epsilon + \gamma \sum_{s'} p(s'|z, b) V^*(s') - p(s'|z, a) V^*(s')$$

the loss at z is

$$\begin{aligned} L_{\tilde{V}}(z) &= V^*(z) - V^{\pi_{\tilde{V}}}(z) \\ (*) &= r(z, a) - r(z, b) + \gamma \sum_{s'} p(s'|z, a) V^*(s') - p(s'|z, b) V^{\pi_{\tilde{V}}}(s') \end{aligned}$$

substitute (*) into (**)

$$L_{\tilde{V}}(z) \leq 2\gamma \epsilon + \gamma \sum_{s'} p(s'|z, b) V^*(s') - p(s'|z, b) V^{\pi_{\tilde{V}}}(s')$$

(the a terms cancel)

$$\begin{aligned} &= 2\gamma \epsilon + \gamma \sum_{s'} p(s'|z, b) [V^*(s') - V^{\pi_{\tilde{V}}}(s')] \\ &= 2\gamma \epsilon + \gamma \sum_{s'} p(s'|z, b) L_{\tilde{V}}(s') \quad \text{defn loss} \\ &\leq 2\gamma \epsilon + \gamma \sum_{s'} p(s'|z, b) L_{\tilde{V}}(z) \quad \text{max loss} \geq \text{other loss} \\ &= 2\gamma \epsilon + \gamma L_{\tilde{V}}(z) \quad \text{put } L_{\tilde{V}}(z) \text{ out, indep of } s' \end{aligned}$$

$$L_{\tilde{V}}(z) \leq 2\gamma \epsilon / (1-\gamma) \quad \text{re-arrange} \quad \square$$

- if we only have data, how get \tilde{V} ? No model parameters to do planning
- how could we estimate V of behavior policy (π_b) used to generate D ?

Monte Carlo estimation, model free, uses complete episodes.
e.g. just look at actual (discounted) sum of returns from trajectories

$$\begin{aligned} s_{i,1}, a_{i,1}, r_{i,1}, s_{i,2}, \dots, r_{i,T} &\rightarrow \sum_{k=1}^T \gamma^{k-1} r_{i,k} \\ s_{i,1,1}, a_{i,1,1}, r_{i,1,1}, s_{i,1,2}, \dots, r_{i,1,T} &\rightarrow \sum_{k=1}^T \gamma^{k-1} r_{i,1,k} \end{aligned}$$

average over return from each traj is an estimation
of $V^{\pi_b}(s_i)$ (assuming all start from same starting s_i)
(empirical mean return instead of expected return)
(empirical mean converges to expectation)

one justification: average converges to expectation

• Chernoff-Hoeffding inequality

Let X_1, X_n be indep random var
assume X_i a.s. bounded, $P(X_i \in [a_i, b_i]) = 1 \quad 1 \leq i \leq n$
then $P(|\frac{1}{n} \sum_{i=1}^n X_i - E[X]| \geq \epsilon) \leq 2 \exp\left(\frac{-2n^2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$

• every visit MC evaluation: everytime s visited in an episode, incremental

• but using MC estimation above only yields estimate of V^{π_b}

what if we want to know $Q^*/V^*/\pi^*$?

need off policy learning

learn about one policy π by potentially executing another
(also relevant to value function)

so can have
finite sample
guarantees on
estimate
quality

to do this, 1st briefly review temporal difference learning

TD methods

model free

learn from incomplete trajectories by bootstrapping

updating V
involves an estimate

MC vs TD

incremental every visit MC. G_t = actual episode return starting at s_t

$$V(s_t) = V(s_t) + \alpha (G_t - V(s_t))$$

TD(0): update $V(s_t)$ towards estimated return $r_{t+1} + \gamma V(s_{t+1})$

$$V(s_t) = V(s_t) + \alpha (\underbrace{r_{t+1} + \gamma V(s_{t+1})}_{\text{TD target}} - \underbrace{V(s_t)}_{\text{TD error}})$$

TD target
TD error

both TD & MC approximate expectation of future rewards
by averaging over samples

in online learning TD has some benefit over MC because don't have to wait until the end of the episode to learn

bias/variance

$G_t = r_{t+1} + \gamma r_{t+2} + \dots$ = unbiased estimate of $V^\pi(s_t)$

true TD target $r_{t+1} + \gamma V^\pi(s_{t+1})$ is an unbiased estimate of $V^\pi(s_t)$

TD target $r_{t+1} + \gamma V(s_{t+1})$ is a biased estimate

because $V(s_{t+1}) \approx V^\pi(s_t)$ (and may have large error w/little data or early in computation)

but G_t / return is much higher variance than TD target

MC high var, 0 bias

TD(0) low var, some bias

in batch data setting for discrete state & actions

MC converges to soln w/min MSE (best fit to observed returns)

$$\sum_{d=1}^D \sum_{t=1}^{H_d} (G_t^d - V(s_t^d))^2$$

TD(0) converges to max likelihood Markov model

[*] for discrete S-A

TD exploits Markov property

MC does not assume nor exploit Markov property } \rightarrow ask about

• Q-learning

like TD(0) but estimates state-action value instead of state value
(useful for extracting a policy)

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha (r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$

converges to optimal Q^* [!] target

o approximate dynamic programming

2 challenges / sources of error

- 1) in our assumed setting, state space S is large or as intractable or impossible in general to represent V exactly
 \Rightarrow instead use function (parametric or nonparametric) give ex of each
 to approximate the Q and/or V function
 yields approximation error - why?
- 2) estimation error. If reward model and dynamics A transition
 model are unknown, cannot exactly compute \rightarrow Bellman backup (an equivalent way to say this is that cannot exactly compute the Bellman operators γ and γ^π) see earlier def
 \Rightarrow instead estimate the Bellman operator from samples

o Popular approach: Fitted Q Iteration (FQI) why useful to compute Q instead of V ?

Gordon 1999: fitted value iteration

Ormoneit & Sen 2002

Ernest et al. 2005 FQI

key idea: pose Q -function determination problem as a sequence of regression problems

recall value iteration in tabular domains

$$Q_{k+1}(s, a) = r(s, a) + \gamma \sum s' p(s'|s, a) \max_{a'} Q_k(s', a')$$

$$Q_k(s, a) = r(s, a)$$

also can express backup as γQ_k where γ is Bellman operator

two ideas. Assume start w/some rep of Q_k

- 1) Approximate backup / expectation using samples & prior
 estimate of Q_k
- 2) fit function approximator (representing Q_{k+1}) to those samples

note: we can get these from it's j

FQI alg
 input: set of N tuples (s_i, a_i, r_i, s'_i) F and a regression algorithm
 where $(s_i, a_i) \sim p$ and $s' \sim p(\cdot | s_i, a_i)$ where p sampling distrib over $S \times A$

where \hat{Q}_k to be a func equal to 0 everywhere on $S \times A$

$k=0$, define \hat{Q}_k to be a func equal to 0 everywhere on $S \times A$

while stopping condition not reached

- $k = k + 1$

- construct new dataset $D'_k = \{(s_i, a_i), y_i\}_{i=1}^N$ where
 $y_i = r_i + \gamma \max_{a'} \hat{Q}_{k-1}(s'_i, a')$ can view as constructing bootstrapped samples of Q
 $= \gamma \hat{Q}_k(s_i, a_i)$

- Use regression algorithm to use D'_k to fit \hat{Q}_k
 $\hat{Q}_k = \text{fit}(f, F, D'_k)$ (find best regression func $f \in F$ that fits data)