

# Problem Setting

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

$S_{i1}, a_{i1}, r_{i1}, S_{i2}, a_{i2}, r_{i2}, \dots$   
 $\uparrow$  time  $\uparrow$  index into traj

gen by behavior  $\pi_b$   
 (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?

Why batch

- 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  $S$ ?  $\swarrow$  describe realistic
- 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  $S, A$ )

thm: (Satinder & Yee 1994). Let  $V^*$  be the optimal value func for a discrete-time MDP having finite state and action sets and an  $\infty$  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  $\forall s \in S$

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

relate to  $Q$

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

proof:

for shorthand define loss  $L_{\tilde{V}}(s) = V^*(s) - V^{\pi_{\tilde{V}}}(s)$   
 $\exists$  a 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  $\pi$  for  $\tilde{V}$ ,  $b$  must look at least as good as  $a$   
 $r(z, a) + \gamma \sum_{s' \in S} p(s'|z, a) \tilde{V}(s') \leq r(z, b) + \gamma \sum_{s' \in S} p(s'|z, b) \tilde{V}(s')$

proof cont

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

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

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

the loss at  $z$  is

$$L_{\tilde{V}}(z) = V^*(z) - V^{\pi_{\tilde{V}}}(z)$$

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

substitute (\*) into (\*\*)

$$L_{\tilde{V}}(z) \leq 2\gamma\epsilon + \gamma \sum_{s' \in S} p(s'|z,b) V^*(s') - p(s'|z,b) V^{\pi_{\tilde{V}}}(s') \quad (\text{the } a \text{ terms cancel})$$

$$= 2\gamma\epsilon + \gamma \sum_{s' \in S} p(s'|z,b) [V^*(s') - V^{\pi_{\tilde{V}}}(s')]$$

$$= 2\gamma\epsilon + \gamma \sum_{s' \in S} p(s'|z,b) L_{\tilde{V}}(s') \quad \text{defn loss}$$

$$\leq 2\gamma\epsilon + \gamma \sum_{s' \in S} p(s'|z,b) L_{\tilde{V}}(z) \quad \text{max loss} \geq \text{other } s'$$

$$= 2\gamma\epsilon + \gamma L_{\tilde{V}}(z) \quad \text{pull } L_{\tilde{V}}(z) \text{ out, indep of } s'$$

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

- if we only have data, how get  $\tilde{V}$ ? No model parameters to do planning / dynamic programming
  - how could we estimate  $V$  of behavior policy ( $\pi_b$ ) used to generate  $D$ ? sample to approx expectation
- Monte carlo estimation, model free, uses complete episodes, e.g. just look at actual (discounted) sum of returns from trajectories

$$s_{11}, a_{11}, r_{11}, s_{12}, \dots, r_{1T}$$

$$\rightarrow \sum_{k=1}^T \gamma^{k-1} r_{1k}$$

$$s_{i1}, a_{i1}, r_{i1}, s_{i2}, a_{i2}, r_{i2}, \dots, r_{iT}$$

$$\rightarrow \sum_{k=1}^T \gamma^{k-1} r_{i+1,k}$$

average over return from each traj is an estimation of  $V^{\pi_b}(s_{i1})$  (assuming all  $s_{i1}$  start from same starting  $s$ )  
(empirical mean return instead of expected return)

one justification: average converges to expectation

Chernoff-Hoeffding inequality

let  $X_1, \dots, X_n$  be indep random var

assume  $X_i$  a.s. bounded,  $P(X_i \in [a_i, b_i]) = 1$

$$\text{then } P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq \epsilon\right) \leq 2 \exp\left(\frac{-2n^2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

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

but using MC estimation above only yields estimate of  $V^{\pi_b}$

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

need off policy learning

learn about one policy  $\pi$  by potentially executing another (also relevant to value function)

to do this, 1st briefly review temporal difference learning (TD)

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  $r_{t+1} + \gamma V(s_{t+1})$

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

TD error

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  $V(s_{t+1}) \neq V^\pi(s_{t+1})$  (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 } → 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  $Q^*$

o approximate dynamic programming

2 challenges / sources of error

1) in our assumed setting, state space  $S$  is large or  $\infty$   
intractable or impossible in general to represent  $V$  exactly  
 $\Rightarrow$  instead use function (parametric or nonparametric)  $\leftarrow$  give ex of each  
to approximate the  $Q$  and/or  $V$  function  
yields approximation error - why?

2) estimation error. If reward model and dynamics / transition model are unknown, cannot exactly compute a Bellman backup (an equivalent way to say this is that cannot exactly compute the Bellman operators  $T$  and  $T^\pi$ )  $\leftarrow$  see earlier def  
 $\Rightarrow$  instead estimate the Bellman operator from samples

o Popular approach: Fitted  $Q$  Iteration (FQI)  $\leftarrow$  why useful to compute  $Q$  instead of  $V$ ?  
Gordon 1999: Fitted value iteration  
Ormoneit & Sen 2002  
Ernst et al. 2005 FQI

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

recall value iteration in tabular domains

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

$$Q_1(s, a) = r(s, a)$$

also can express backup as  $TQ_k$  where  $T$  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

$\leftarrow$  note: we can get these from 1)

FQI alg

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

$k=0$ , 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

$$y_i = r_i + \gamma \max_{a'} \hat{Q}_{k-1}(s'_i, a') \quad \left\{ \begin{array}{l} \text{can view as constructing} \\ \text{bootstrapped samples of } Q \end{array} \right.$$
$$\equiv \hat{T} \hat{Q}_k(s_i, a_i)$$

- Use regression algorithm to use  $D'_k$  to fit  $\hat{Q}_k$

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

essentially have a

target function  $= T Q_k$

noisy observ  $= \hat{T} Q_k$  at a finite set of points

minimize empirical error in the fitting process

$$\hat{Q}_k = \min_{f \in \mathcal{F}} \|f - \hat{T} \hat{Q}_{k-1}\|_{\hat{p}}, \quad \hat{p} = \text{empirical measure of } p \text{ given } N \text{ samples}$$

regressor class

with the goal of minimizing the true error

$$Q_k = \min_{f \in \mathcal{F}} \|f - T Q_{k-1}\|_p$$

regression problem of  $\hat{Q}_k = T \hat{Q}_{k-1}$

objective is  $\|\hat{Q}_k - T \hat{Q}_{k-1}\| = \underbrace{\|\hat{Q}_k - Q_k\|_p}_{\text{estimation error}} + \underbrace{\|Q_k - T \hat{Q}_{k-1}\|_p}_{\text{approximation error}}$  \* important idea

doesn't address distance to true  $Q$  w/ no func approx

- what should the error function be?
- what should the regression function be?
- does it matter

for desirable properties

for performance

- what if  $p$  (sampling distrib) is biased? how could this affect things?

example: least squares value iteration (from Tsitsiklis & Van Roy 1996)

choose sum of squared errors as error function

compute parameters to find least squares fit

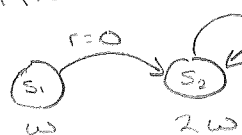
let  $w$  be the parameters that define a function  $f$

then  $\hat{Q}_k = \text{fit}(f_i, F_i, D'_k)$  becomes

$$w_k = \argmin_w \sum_{i=1}^n (\hat{Q}_k, w_k(s_i, a_i) - \hat{T} \hat{Q}_{k-1, w_{k-1}}(s_i, a_i))^2$$

how good is least-squares value/q iteration?

ex. (from Tsitsiklis & Van Roy)



2 states  
1 action

all rewards 0

$$p(s_2 | s_1, a_1) = 1$$

$$p(s_2 | s_2, a_1) = 1$$

$$Q^*(s_1, a_1) = V^*(s_1) = V^*(s_2) = Q^*(s_2, a_1) = 0$$

define a feature  $\phi$  over state space s.t.  $\phi(s_1) = 1$   $\phi(s_2) = 2$

use a linear function approximator of the (state/action) value

$$\hat{Q}_{k, w_k}(s) = w_k \phi(s) \quad \text{where } w_k \text{ is a scalar}$$

1) can we represent optimal value func w/ this approx?

2) assume a discount factor  $\gamma$

imagine 2 data tuples:  $s_1, a_1, 0, s_2$   
 $s_2, a_1, 0, s_2$

if simple  
linear approx  
works

$$w_{k+1} = \argmin_w \sum_{i=1}^2 (w \phi(s_i) - [\gamma w \phi](s_i))^2$$

define  $T w \phi$

solve for  $w$

3) how does  $w$  change as  $k$  gets larger?

$$\begin{aligned}
 \omega_{k+1} &= \operatorname{argmin}_{\omega} \sum_{i=1}^2 (\omega \phi(s_i) - [\mathcal{T}\omega\phi](s_i))^2 \\
 &= (\omega - (0 + \gamma 2\omega_k))^2 + (2\omega - (0 + \gamma 2\omega_k))^2 \\
 &= \omega^2 - 4\omega_k\gamma\omega + \gamma^2 4\omega_k^2 + 4\omega^2 - 8\omega\omega_k\gamma + 4\gamma^2\omega_k^2
 \end{aligned}$$

$$d/d\omega \rightarrow 10\omega - 4\omega_k\gamma - 8\omega_k\gamma = 0$$

$$10\omega = 12\omega_k\gamma$$

$$\omega = 6/5\omega_k\gamma$$

$$w_{k+1} = \arg \min_w \left[ \underbrace{(w - \gamma^2 w_k)^2}_{\text{from } S_1} + \underbrace{(2w - \gamma^2 w_k)^2}_{\text{from } S_2} \right]$$

$d/dw$  & set to 0

$$w_{k+1} = 6/5 \gamma^2 w_k$$

if  $\gamma > 5/6$  then  $w_{k+1} \rightarrow \infty$  if  $w \neq 0$

diverges!

Why does this occur?

Bellman backup is a contraction

but projecting back to func approximator class may not be

see Geoff Gordon

what function approximators do guarantee convergence of FQI?  
averagers / kernels / nearest neighbors

even if not guaranteed to converge, can do well

guarantees if can bound error per iteration of FVI (similar to FQI) Mohammad

assume we can approximate Bellman optimality operator up to  $\epsilon$  at iteration  $k$

$$\hat{V}_{k+1} = \hat{T} V_k \quad \text{s.t.} \quad \hat{V}_{k+1} = T V_k + \epsilon_k \quad \text{or} \quad \|\hat{V}_{k+1} - T V_k\|_\infty = \epsilon_k$$

→ Mohammad

also exist work on finite sample bounds

#### • summary

- FVI / FQI popular and powerful algorithm
- performing/solving supervised learning problem at each iteration
- due to iterations, error can compound (worse than linear) across iterations
- input data sampling distrib  $p$  has a crucial influence on algorithm performance. (equivalent to role of exploration)
- poor performance (divergence) is possible: choice of error function & regressor important
  - issue also of capacity of regressor vs true function vs available data - to be discussed more!
- falls into class of model free value function based techniques