

03/11/14

Review: (alg. b/w off constraint prob not known)

\* Markov decision process

$$MDP = \{S, A, P(s'|s, a), R(s), \gamma\}$$

\* Value functions

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right] \quad (\text{Expected value})$$

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

\* Bellman Optimality Equation:

$$V^*(s) = R(s) + \gamma \max_a \sum_{s'} p(s'|s, a) V^*(s')$$

\* Algorithm: Value iteration ← for Computing  $V^*(s)$  directly.

$$\rightarrow \text{initialize } V_0(s) = 0 \quad \forall s$$

$$\rightarrow \text{iterate } V_{k+1}(s) = R(s) + \gamma \max_a \left[ \sum_{s'} p(s'|s, a) V_k(s') \right]$$

To prove: The algorithm converges

$\rightarrow V^*(s)$  is a fixed point of ~~last~~ iteration

$\rightarrow$  There are no other ~~fixed~~ points.

$\rightarrow$  The algorithm always reaches  $V^*(s)$

Before going to the proof of Converges, we will prove a lemma

$$(a)^* v \leq \{ (a)_i v \}_{i \in A} \quad \forall v \in \mathbb{R}^A$$

Lemma: for any functions  $f(a)$  and  $g(a)$

$$|\max_a f(a) - \max_{a'} g(a')| \leq \max_a |f(a) - g(a)|$$

Proof of lemma:

$$f(a) - \max_{a'} g(a') \leq f(a) - g(a), \text{ for all } a$$

$$\max_a f(a) - \max_{a'} g(a') \leq \max_a (f(a) - g(a)), \text{ max over } a$$
$$\leq \max_a |f(a) - g(a)| \rightarrow \textcircled{1}$$

By symmetry (exchanging  $f$  &  $g$ )  $\rightarrow \textcircled{1} \Rightarrow \textcircled{2}$

$$\max_a g(a) - \max_a f(a) \leq \max_a |f(g(a)) - f(a)| \rightarrow \textcircled{2}$$

↳ same as  $|f(a) - g(a)|$

The L.H.S. of eq $\textcircled{1}$  or eq $\textcircled{2}$ , must correspond to absolute value of

$$\max_a g-f(a) - \max_a g(a)$$

thus, we get

$$|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$$

Using this lemma, we will prove

Theorem: Value Iteration Converges.

$$\Rightarrow \lim_{k \rightarrow \infty} [v_k(s)] \rightarrow v^*(s) \text{ for all states, } s$$

Proof: let  $\Delta_k = \max_s |v_k(s) - v^*(s)|$

$$\Delta_{k+1} = \max_s |v_{k+1}(s) - v^*(s)|$$

$$= \max_s \left| R(s) + \gamma \sum_{s'} P(s'|s,a) v_k(s') - v^*(s) \right|$$

$$= \gamma \max_s \left| \max_a \sum_{s'} P(s'|s,a) v_k(s') - \max_a \sum_{s'} P(s'|s,a) v^*(s') \right|$$

using the lemma

$$\leq \gamma \max_s \max_a \left| \sum_{s'} P(s'|s,a) [v_k(s') - v^*(s')] \right|$$

$$\leq \gamma \max_s \max_a \left| \left( \sum_{s'} P(s'|s,a) \right) \max_{s''} |v_k(s') - v^*(s'')| \right|$$

$$= \gamma \max_s \max_a |\Delta_k|$$

$$= \gamma \Delta_k \quad [\Delta_k \text{ does not depend on } a \text{ or } s]$$

By induction, we can show

$$\boxed{\Delta_k \leq \gamma^k \Delta_0}$$

Thus  $\Delta_k \rightarrow 0$  as  $k \rightarrow \infty$

Assume all the rewards are bounded.

$$\Delta_0 = \max_{s \in S} |\nu_0(s) - \nu^*(s)| = \max_{s \in S} |\nu^*(s)| \quad [\because \nu_0(s) = 0]$$

$$\leq \max_{s \in S} |R(s)| \left[ 1 + r + r^2 + \dots \right]$$

$$= \max_{s \in S} |R(s)| \left( \frac{1}{1-r} \right)$$

Thus  $\boxed{\Delta_k \leq \frac{1}{1-r} \max_{s \in S} |R(s)|}$   $\rightarrow 0$  as  $k \rightarrow \infty$

The above inequality suggests that the convergence rate

$\rightarrow$  depends on  $r$

$\rightarrow$  more iterations are required as  $r \rightarrow 1$

FINAL COVERS UP TO HERE !!

### REINFORCEMENT LEARNING (not for final)

\* What if  $P(s'|s, a)$  and  $R(s)$  are not known?

Can we learn  $\pi^*(s)$  or  $\nu^*(s)$  from experience?

With no reward function

$$\boxed{\Delta^* = \Delta}$$

and  $\Delta$  is defined as

- 1) Model-based (indirect) approach
- Explore world, estimate model  $\hat{P}(s'|s,a) \approx P(s'|s,a)$ ,  
 Compute  $\hat{\pi}^*(s)$  or  $\hat{v}^*(s)$  from  $\hat{P}(s'|s,a)$
- \* Cons:  $\rightarrow$  to store  $P(s'|s,a)$  is  $O(n^2)$  for n. states  
 $\rightarrow$  only care about  $\pi^*(s)$  or  $v^*(s)$  which are  $O(n)$   
 $\rightarrow$  Is it really necessary to estimate a model?

\* Pros:  $\rightarrow$  model  $P(s'|s,a)$  useful for task transfer,  
 settings whose rewards  $R(s)$  or discount factor  $r$  changes,  
 but  $P(s'|s,a)$  stay the same.

- 2) Direct Approach: learn  $\pi^*(s)$ ,  $v^*(s)$  without building a model

### Stochastic Approximation Theory

\* To estimate mean of random variable,  $x$  from samples  $x_0, x_1, x_2, \dots, x_{T-1}$

1) obvious  $\rightarrow$  sample average

$$\mu = \frac{1}{T} [x_0 + x_1 + \dots + x_{T-1}] \rightarrow E[x] \text{ (as } T \rightarrow \infty)$$

Converges by law of large numbers

## 2) Incremental update:

→ Initialize  $\mu_0 = 0$

→ update  $\mu_t = (1-\alpha_t) \mu_{t-1} + \alpha_t r + x_t$  for  $0 < \alpha_t < 1$

obtain  $\mu$  not  $(\mu_t)$  or  $(\mu_t)$  more compact

This can be also written as

$$\mu_t = \mu_{t-1} + \alpha_t (\hat{x}_t - \mu_{t-1})$$

temporal difference

intuitively this algorithm is called Temporal Difference algorithm.

Theorem:  $\mu \rightarrow E[x]$  as  $t \rightarrow \infty$  if

i)  $\alpha_t$  doesn't decay too fast  $\rightarrow \sum_{t=1}^{\infty} \alpha_t = \infty$

ii)  $\alpha_t$  doesn't decay too slow  $\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 < \infty$  is finite

One possible choice of  $\alpha_t$ :

choose  $\alpha_t = \frac{1}{t}$ . ~~also~~ satisfies the conditions.

This actually gives the same estimate as the running average.

and so