

# Reinforcement Learning (Learning Algos)

Friday, March 27, 2020 3:05 PM

Previously, we covered the planning problem (aka "Dynamic Programming")

which finds a good policy  $\pi$  assuming the environment (reward + transition probabilities) is known.

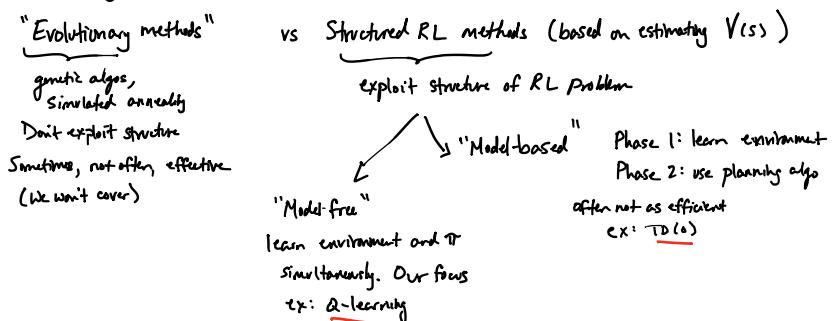
This was "tractable" (eg, equiv. to a LP) though still have difficulties when  $|S| = 10^{400}$  ...

Now, consider learning problems, where not only do we want  $\pi$ , but the environment is unknown.

Simultaneously learn  $\pi$  and transition probabilities. There is now a clear exploration-exploitation tradeoff.

This is fundamentally harder.

Types of learning methods:



Background: Stochastic Approx.

Law of Large Numbers (LLN)

Let  $X_i$  be iid <sup>or mutually independent</sup> realizations of a r.v.  $X$  s.t.  $E[X] = \mu$  and  $E[X^2] < \infty$

Define  $\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m X_i$  (note:  $\hat{\mu}_{m+1} = \underbrace{\frac{1}{m+1} X_{m+1}}_{\alpha} + \underbrace{\frac{m}{m+1} \hat{\mu}_m}_{1-\alpha}$ )  
then

Weak LLN:  $\hat{\mu}_m \xrightarrow{P} \mu$ , i.e.,  $\lim_{m \rightarrow \infty} P(|\hat{\mu}_m - \mu| > \varepsilon) = 0$  ( $\forall \varepsilon > 0$ )

Strong LLN:  $\hat{\mu}_m \xrightarrow{as.} \mu$ , i.e.,  $P(\lim \hat{\mu}_m = \mu) = 1$

Theorem 17.4 (mean estimation, generalization of SLLN)

$X \in [0, 1]$  a r.v., w/ iid realizations  $\{X_0, \dots, X_m\}$ , and define  $\hat{\mu}_{m+1} = (1 - \alpha_m) \hat{\mu}_m + \alpha_m X_m$   
 $\Rightarrow \hat{\mu}_0 = X_0$ ,  $\alpha_m \in [0, 1]$ ,  $\sum \alpha_m = +\infty$ ,  $\sum \alpha_m^2 < \infty$  (ex:  $\alpha_m = \frac{1}{m}$ )  
 then  $\hat{\mu}_m \xrightarrow{as.} E[X]$  ( $\Rightarrow \hat{\mu}_m \xrightarrow{L^2} E[X]$  quadratic mean, aka  $L^2$ , convergence means)  
 $\lim_{m \rightarrow \infty} E[(\hat{\mu}_m - E[X])^2] = 0$

proof (of  $g_m$ , not  $a.s.$ , for simplicity)

- $E[\mu_n] = E[X_n] = E[X]$ , and by linearity and induction,  $E[\mu_m] = E[X]$  also thus convergence in  $L^2$  just means  $\lim_m \text{Var}[\mu_m] = 0$
  - By independence,  $\text{Var}(\mu_{m+1}) = (1-d_m)^2 \text{Var}(\mu_m) + d_m^2 \cdot \underbrace{\text{Var}(X)}_{=1}$   
 $\leq (1-d_m)^2 \text{Var}(\mu_m) + d_m^2 \leq \underbrace{(1-d_m)}_{\text{so } t^2 \leq t} \text{Var}(\mu_m) + d_m^2$

We're going to start by showing  $\liminf_m \text{Var}(\mu_m) = 0$

... Suppose not. Since  $\text{Var}(\cdot) \geq 0$ , this means  $\exists \varepsilon > 0, \exists N \text{ s.t. } \forall m > N, \text{Var}(\mu_m) \geq \varepsilon$

Wlog, relabel sequence so  $N=0$

$$\begin{aligned}
 \text{Var}(\mu_{m+1}) &= \text{Var}(\mu_m) - d_m \underbrace{\text{Var}(\mu_m)}_{\geq 0} + d_m^2 \\
 &\leq \text{Var}(\mu_m) - d_m^2 + d_m^2 \\
 &\leq \text{Var}(\mu_{m-1}) - \varepsilon(d_m + d_{m-1}) + d_m^2 + d_{m-1}^2 \quad (\text{recurred}) \\
 &\dots \\
 &\leq \text{Var}(\mu_0) - \varepsilon \sum_{k=0}^m d_k + \sum_{k=0}^m d_k^2 \\
 \Rightarrow \lim \text{Var}(\mu_{m+1}) &= -\infty, \text{ impossible.} \quad \begin{matrix} \nearrow +\infty \\ \text{bounded} \end{matrix}
 \end{aligned}$$

Impression

Now, since  $\sum \alpha_m^2 < \infty \Rightarrow \alpha_m \rightarrow 0$ , so  $\exists M$  s.t.  $\alpha_m < \varepsilon + m/M$

So,  $\exists m > M$  s.t.  $\text{Var}(\mu_m) < \varepsilon$  ( $\not\equiv \forall m > M \dots$  that's what we need to show)

$$\text{Again, use } \text{Var}(\mu_{m+1}) \leq (1-\alpha_m) \underbrace{\text{Var}(\mu_m)}_{\leq \varepsilon} + \alpha_m^2 \Rightarrow \alpha_m^2 = \alpha_m \cdot \alpha_m < \varepsilon \cdot \alpha_m \\ \leq (1-\alpha_m)\varepsilon + \alpha_m\varepsilon = \varepsilon$$

$\Rightarrow \forall m \geq M, \text{Var}(\mu_m) \leq \varepsilon$  □

## Stochastic Approximation

Solves things like  $x = H(x)$  where  $H(x)$  can't be queried directly, but we can

observe noisy measurements,  $y_i = H(x_i) + w_i$ ,  $E[w_i] = 0$

Ex:  $\min f(x) := Ef(x, w)$      $w$  over  $\omega$

If  $f$  is convex and smooth, minimizer is characterized by  $0 = \nabla f(x)$  case

i.e.,  $x = \underbrace{(I - P_f)}_H(x)$ . Under some conditions (ex, discrete distribution, or,  $P_f(x, \omega)$  is uniformly Lipschitz cts w.r.t  $\omega$ ) then  $P_f(x) := P F_f(x, \omega) = E P_f(x, \omega)$

So if we can sample  $w$ , we can sample  $\nabla f$  and hence  $H$ .

so, iteration

$$\text{So, iteration } X_{t+1} = (1-\alpha_t)X_t + \alpha_t(H(X_t) + w_t)$$

looks like Thm 1F.4

## Background:

Martingale (discrete-time, but can be defined for cts. time also)

$(X_t)_{t \in \mathbb{N}}$  a sequence of r.v. w/  $\mathbb{E}[|X_t|] < \infty$

$$\text{st. } \mathbb{E}[X_{t+1} | X_t, X_{t-1}, \dots, X_0] = X_t \quad \text{Ex. a random walk}$$

" "  $\equiv x_n$  Supermarktgäste

$$\begin{array}{ll}
 " & " \geq X_n \quad \text{Submartingale} \\
 " & " = 0 \quad \text{Martingale difference} \\
 & \text{ex: if } (X_n) \text{ is a Martingale,} \\
 & Y_n = X_n - X_{n-1} \text{ is a Martingale Difference Sequence}
 \end{array}$$

Fact (Thm 17.5, supermartingale convergence)

$(X_t)_{t \in \mathbb{N}}, (Y_t)_{t \in \mathbb{N}}, (Z_t)_{t \in \mathbb{N}}$  nonneg. r.v. &  $\sum Y_t < \infty$ , and  $\mathcal{F}_t$  the filtration generated by all r.v. for  $t' \leq t$

Then if  $E[X_{t+1} | \mathcal{F}_t] \leq X_t + Y_t - Z_t$ , then

- ①  $X_t$  converges w.p.1 (i.e.,  $Y_t, Z_t \rightarrow 0 \Rightarrow$  supermartingales converge)
- ②  $\sum Z_t < \infty$

Use this to show Thm 17.16, 17.17 which generalizes our simple SA "LLN"-like theorem  
... too far off topic, so see book (Mahajan et al.)

### Temporal Difference algorithm (TD(0))

most key (and novel) idea in RL, combining dynamic programming ideas w/ Monte Carlo (Sutton + Barto)

All based on estimating  $V(s)$

( $V$  is a function of the policy  $\pi$ ,  $V_\pi(s)$ , but if we find  $V'(s)$  ( $= V_{\pi'}(s)$ )

then we can recover  $\pi'$ , like in the greedy update of Policy Iteration,  $\pi' \leftarrow \arg \max_\pi R_\pi + \gamma P_\pi V$ )

*his linear ones*

Based on Eq 17.5 / 17.6 ("Bellman's Eq.", but not his optimality one)

matrix notation

$$\begin{aligned}
 V_\pi(s) &= E[r(s, \pi(s))] + \gamma \sum_{s' \in S} P[s'|s, \pi(s)] \cdot V_\pi(s') \quad \text{i.e., } V_\pi = R_\pi + \gamma P_\pi V_\pi \\
 &= E_{s'}[r(s, \pi(s)) + \gamma \cdot V_\pi(s') | s] \quad \therefore V_\pi = (I - \gamma P_\pi)^{-1} R_\pi
 \end{aligned}$$

... but  $s' \sim P[s'|s, \pi(s)]$  is unknown, but we can sample from it

[setting: a model of the environment is unknown, but we still interact w/ it]

We'll start w/ a guess for  $V$  (so  $V(s)$  defined  $\forall s \in S$ ), assume deterministic policies

- then sample a new state  $s'$  (environment does this for us, as a black box)
- reward  $r(s, \pi(s))$  is also sampled (by the environment)
- we update policy value

$$\begin{aligned}
 V(s) &\leftarrow (1-\alpha) V(s) + \alpha(r(s, \pi(s)) + \gamma V(s')) \quad \xrightarrow{\text{Bellman's Linear Eqn}} E(r(s, \pi(s)) + \gamma V(s')) = V(s) \\
 &= V(s) + \alpha \underbrace{(r(s, \pi(s)) + \gamma V(s') - V(s))}_{\text{temporal difference of } V \text{ values}} \quad \cdots \text{so this should get small} \\
 &\quad \alpha \text{ depends on} \quad \text{as } V \text{ is close to being accurate} \\
 &\quad \text{how many times} \\
 &\quad \text{we've visited state } s \text{ before}
 \end{aligned}$$

Algo: TD(0) Input:  $\pi$

Initialize  $V$  arbitrary, " $t$ " is now an epoch, like a trial or a game

For  $t = 0, 1, \dots, T$

$S$  initiated arbitrarily/randomly (new game)  
(other than a terminal state)

for each time step within epoch  
 (rounds of the game) eg robot simulation  
 observe ("sample")  $r' = \text{Reward}(s, \pi(s))$  from environment  
 observe ("sample")  $s'$  as our next state, from environment  
 $\hat{V}(s) \leftarrow (1-\alpha)\hat{V}(s) + \alpha(r' + \gamma\hat{V}(s'))$   
 $s \leftarrow s'$   
(stop when reach a terminal state)  
 end      end

Finite-time MDP has a final state  
 (ex: in a board game, winning/losing is final state)  
 may not know a priori how many rounds it will last  
 For  $\infty$ -horizon, maybe approximate?

You can prove convergence via Thm. 17.17  
 $\rightarrow (\forall \pi), V(\pi) \rightarrow V_\pi(s) \quad \forall s \in S$   
 ... but how to choose  $\pi$  not clear.  
both taken care of by Q-learning framework.

### Q-learning

If we know <sup>optimal</sup> state-action value function  $\hat{Q}^*(s, a)$ , then  
 1) it gives  $V^*(s) = \max_{a \in A} Q^*(s, a)$   
 2) it gives  $\pi^*(s) = \arg \max_{a \in A} Q^*(s, a)$

For simplicity, now assume deterministic reward function.  
 Bellman optimality equations, (Eq (17.4) + Eq (17.1))  

$$\begin{aligned} Q^*(s, a) &= \mathbb{E}[r(s, a) + \gamma \sum_{s' \in S} P[s'|s, a] V^*(s')] \\ &= \mathbb{E}[r(s, a) + \gamma \max_{a'} Q^*(s', a')] \xleftarrow{*} (*) \end{aligned}$$

### ALGO: Q-learning

Initialize  $Q$  arbitrarily // no  $\pi$  needed

For epochs  $t = 0, 1, \dots, T$

Pick a state  $s$  // same as TD(0)

For each round in the epoch

Choose action  $a$  based on  $\pi, s$

Observe reward  $r'$  from environment

Observe state  $s'$  from environment

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

end when reach a terminal state

End

End

So, a bit like TD( $\alpha$ ) but update  $Q$ , not  $V$ , and  $\pi$  doesn't need to be pre-specified

Update can also be written  $Q(s, a) \leftarrow (1-\alpha)Q(s, a) + \alpha(r(s, a) + \max_{a'} Q(s', a'))$

① update rule      old      cf (\*) Stochastic Approx.

### ③ learning policy

#### Policy $\pi$

- For general convergence, this is arbitrary as long as every pair  $(s, a)$  is visited  $\infty$ -ly often

(... of course, at convergence, recover  $\pi^*$  via maximization in 2) above )

$$\pi^*(s) \in \arg\max_a Q^*(s, a)$$

- For efficient convergence,

$$\pi(s) \in \arg\max_a Q_t(s, a) \text{ is natural choice ...}$$

... but this is all **exploitation** and if

$Q_t$  isn't optimal, we may be missing important  $(s, a)$  pairs

so

" $\epsilon$ -greedy policy" is to take this greedy step w.p.  $1-\epsilon$ ,

and take a random step  $s$  (e.g., uniformly at random) w.p.  $\epsilon$   
exploration

or ~ Boltzmann exploration ... ("temperature"  $\rightarrow 0$  as  $t \rightarrow \infty$ )

Many other variations on ① update rule and ② learning policy

- TD( $\lambda$ )

- on-policy methods like SARSA

- See Sutton and Barto for a lot more

(for  $\alpha$ -learning)

Thm 17.18 For a finite MDP, if  $\forall s, \forall a$   $\sum_{t=0}^{\infty} \alpha_t(s, a) = +\infty$ ,  $\sum_{t=0}^{\infty} \alpha_t^2(s, a) < \infty$   
 $\alpha_t \in [0, 1]$  then  $Q \rightarrow Q^*$  a.s.

$\uparrow t$  is epoch in algo, not rounds in the game.

$$\text{i.e. } \alpha_t(s, a) = \begin{cases} \alpha_t(s_t, a_t) & \text{if } s=s_t, a=a_t \\ 0 & s \neq s_t \text{ or } a \neq a_t \end{cases}$$

so  $\sum_t \alpha_t(s, a) = \infty \forall (s, a)$   $\Rightarrow$  every  $(s, a)$  visited  $\infty$ -ly often.

Proof sketch: use Thm 17.17. See Mohri

## Large State Spaces

Still a problem, as they were for planning algorithms

Very active research

Example approach, map  $\Phi: S \rightarrow \mathbb{R}^{N \approx 100}$

$$|S| \approx 10^{100}$$

$$V_\pi(s) \approx f_\omega(s),$$

$$f_\omega(s) = \langle \omega, \Phi(s) \rangle \text{ for example.}$$

$$\text{Learn } \omega$$