

The Simplest RL in an MDP

Notation

→ in X_t ; choose A_t

→ get $(X_{t+1}, R_{t+1}) \sim P(\cdot | X_t, A_t)$

↑
New State ↑
Reward

IID

likely not known

Q: estimate

$$V^\pi(x) = E \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid X_0 = x \right]$$

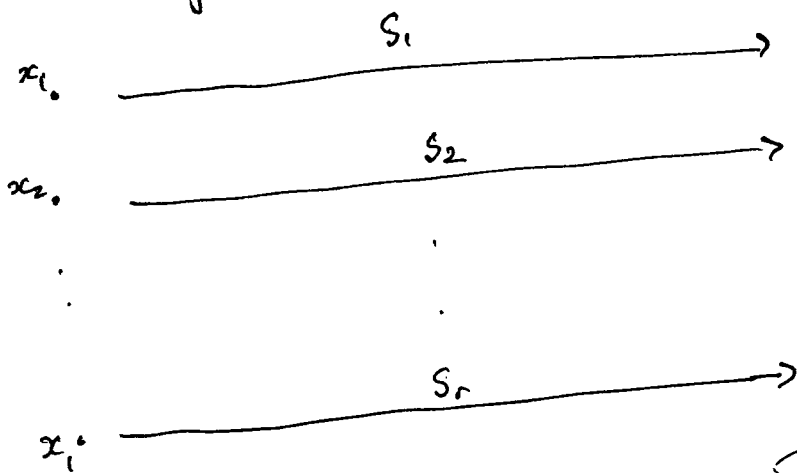
here you act under π .

Simplest

MC estimate

- Start at x , run π , observe reward.
(this assumes MDP has a terminal state)
- Do this many times, average.

and average.



I could form $\frac{1}{\#}$ \sum (Reward of segs starting at x_i)
 → but what if I ~~form~~ get a new seg.?

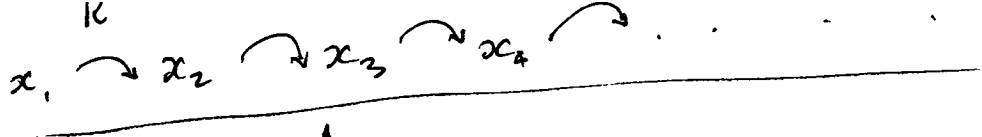
instead

$$\hat{V}^\pi(x) \rightarrow \hat{V}^\pi(x) + \alpha \left[\text{Reward of new seg starting at } x - \hat{V}^\pi(x) \right]$$

(Easy) This is unbiased est.

(Easy) But may not be very efficient

Because we have ests of \hat{V}^π (something in x 's future)



MC $\hat{V}^\pi(x_1)$ est. by sum of discounted rewards along seq. (say) $\hat{V}^\pi(x_2)$

So another est. is $R + \gamma \hat{V}^\pi(x_2)$.
 This has to be unbiased, and might be better

TD(0)

$$S_{t+1} = R_{t+1} + \gamma \hat{V}_t(x_{t+1}) - \hat{V}_t(x_t)$$

$$\hat{V}_{t+1}(x) = \hat{V}_t(x) + \alpha S_{t+1} \mathbb{1}_{[x_t = x]}$$

- This is temporal difference learning.
- ~~is~~ easy to adapt to off-policy seqs.

assume you have

(X_t, R_{t+1}, Y_{t+1}) where

$X_t \sim$ some arb. ergodic MC over X .

$(Y_{t+1}, R_{t+1}) \sim P(\cdot | X_t, A_t)$
↑ from the policy.

then

$S_{t+1} = R_{t+1} + \gamma \hat{V}(Y_{t+1}) - \hat{V}(X_t)$

works.

Which is better?

Depends

example in Szepesvari, p 22
Shows small change in model that
favors either TD(0), MC, depending
on state.

Notice that $TDC(0)$ is a family of algs.

MC: est by counting reward from a start

est by multiple steps of reward,
then pass to est = $TDC(\lambda)$.

$TDC(0)$: est by reward from one step, then add est. of new state

$TDC(\lambda)$ — eqns p 24 of Szepesvári.

What if there are too many states:

eg. continuous, etc.

• estimate $\hat{V}(s)$ with some form of function approximator. (Neural network these days)

TD(λ) then updates θ

rather than V directly
(Sutton & Barto, p 21 gives eqns.)
→ this can diverge

but this is two fairly natural strategies.

- Build a model of the Q function

$$Q^*(x, a) = E_{(\text{best policy})} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid X_0 = x, A_0 = a \right]$$

- particularly nice if I can do this without looking at seqs. observe $X_t, A_t, R_{t+1}, Y_{t+1}$

$$S_{t+1} = R_{t+1} + \gamma \max_{a' \in A} Q(Y_{t+1}, a') - Q(X_t, A_t)$$

$$Q(x, a) \rightarrow Q(x, a) + \alpha S_{t+1} \cdot \mathbb{I}_{\{x=X, a=A\}}$$

NOTICE : - TD update

$$- (Y_{t+1}, R_{t+1}) \sim P(\cdot \mid X_t, A_t)$$

- But we don't need a seq (as long as

X_t is ergodic.

eg. i.e. you can build a local model, which is rather nice for many problems

Alternative: build a model of policy, update.

policy model: parametric probability model $\pi_{\theta} = p(a_t | s_t, \theta)$

common is some form of neural network predicting means, var log sds of a Gaussian policy.

update by gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\pi_{\theta}) \Big|_{\theta_k}$$

↑
expected return

$$J(\pi_\theta) \Big|_{\theta} = \underbrace{E_{z \sim \pi_\theta} [R(z)]}_{\text{expected return under } \pi_\theta}$$

~~assume~~ y_0

write $z = (s_0, a_0, s_1, \dots, s_{T+1})$ trajectory

$$P(z | \theta) = p_0(s_0) \cdot \prod_{t=0}^T \left[P(s_{t+1} | s_t, a_t) \cdot \pi_\theta(a_t | s_t) \right]$$

↑
prior of starting in 0

• Notice $\nabla_{\theta} \log f(\theta) = \frac{1}{f(\theta)} \cdot \nabla_{\theta} f(\theta)$

logarithmic diff reinvented!

$$\log P(\tau|\theta) = \log P_0 + \sum_{t=0}^{\tau} \left\{ \log P(s_{t+1} | a_t, s_t) + \log \pi_{\theta}(a_t | s_t) \right\}$$

\uparrow
NO θ !

So $\nabla_{\theta} \log P(\tau|\theta)$

$$= \sum_{t=0}^{\tau} \left\{ \log \pi_{\theta}(a_t | s_t) \right\}$$

So:

$$\begin{aligned} \nabla_{\theta} J &= \nabla_{\theta} \int_{\tau} P(\tau|\theta) R(\tau) \\ &= \int_{\tau} \nabla_{\theta} P(\tau|\theta) R(\tau) \\ &= \int_{\tau} P \nabla_{\theta} \log P R(\tau) \\ &= E_{\tau} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t | s_t) R(\tau) \right] \end{aligned}$$

Expectation!

But this is
reward for
whole trajectory

some funny thing in approximation
(OpenAI notes) gets

$$\nabla_{\theta} J = E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi \left(\sum_{i=t}^T R(s_i, a_i, s_{i+1}) \right) \right]$$

this is a better estimate.:

we had $\nabla_{\theta} J = E_{\tau} [f_1(\tau) + f_2(\tau)]$

BUT $E_{\tau} [f_1(\tau)] = 0$

→ removing f_1 reduces variance - this is not just manipulation for showing off.

an important strategy.

assume I want to estimate $E_p[h] =$

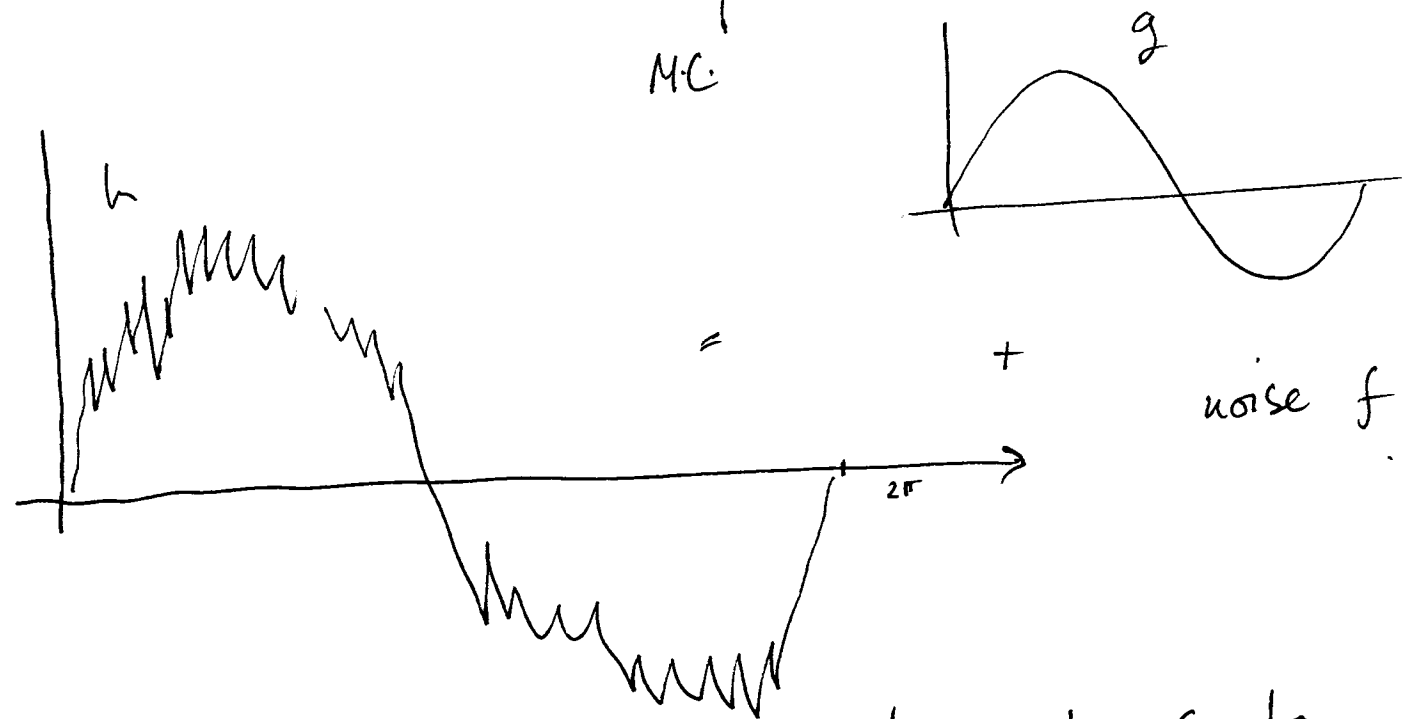
$$E_p [f + g]$$

↑
↑

hard to integrate
easy to integrate.

Always use $E_p [f] + [\text{easy est of } g]$

|-----|
↑
MC



→ notice that h adds large terms to sum, which then cancel (or should, but don't precisely)

Now

$$E_{a_t \sim \pi} \left[\nabla_{\theta} \log \pi(a_t | s_t) \cdot b(s_t) \right]$$

$$= E_{a_t \sim \pi} \left\{ \int \left[P(a_t | s_t) \cdot \nabla_{\theta} \log \pi(a_t | s_t) \right] \cdot b(s_t) \right\}$$

$$\downarrow$$

$$\pi(a_t | s_t)$$

BUT $\int P \nabla_{\theta} \log P = \int \nabla_{\theta} P = \nabla_{\theta} \int P = 0$

$$\therefore E_{a_t \sim \pi} \left[\nabla_{\theta} \log \pi(a_t | s_t) b(s_t) \right] = 0$$

often referred to as a baseline

this means

$$E_{a_t \sim \pi} \left[\nabla_{\theta} \log \pi(a_t | s_t) \cdot b(s_t) \right] = 0$$

Natural choice of measure

$$V_t^\pi$$

est as above

$$\nabla_{\theta} \hat{J}(\pi_{\theta}) = E_{\tau} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi(a_t | s_t) \left(\sum_{u=t}^T R(s_u, a_u, s_{u+1}) - V^{\pi}(s_t) \right) \right]$$

Other base lines

$Q^{\pi}(s_t, a_t)$ works.

$Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$ works

Which is best?

Don't know

- BE CAREFUL
 - the reason there are lots of methods is they mostly DON'T WORK RELIABLY.

• best cases seem to be ones you could do with LQR, but the linearization is a nuisance.

- velocity control of vehicle

- stopping at right place from a single view of stop sign

- important features:
 - local reward is quite important
 - you can generate a ton of examples from simulator
 - view control is easy.

however

You can estimate $\text{Var}[r]$, $\text{Cov}[]$
from MC replicates. to get variance reduction.

This suggests a simple way to get
improved var in Policy gradient, cause
you know \int baseline.

This idea extends:

- under some circumstances, you don't
need to know r
(g must be easy to integrate,
and a covariance cond)

- you can use multiple control
variables.

Control Variates

we want to know

$$I_f = \int f P dx = E_p[f]$$

assume we have g , st $E_p[g] = \gamma$
↑
KNOWN

then

$$\frac{1}{N} \sum_i [f(x_i) + c(g(x_i) - \gamma)] = I_f + \xi$$

where

$$E[\xi] = 0$$

$$\begin{aligned} \text{Var}(\xi) &= \text{Var}[\text{est of } f \text{ from } \frac{1}{N} \sum_i f(x_i)] \\ &+ c^2 \text{Var}[\text{est of } \gamma \text{ from } \frac{1}{N} \sum_i g(x_i)] \\ &+ 2c \text{Cov}[\text{est of } f, \text{est of } \gamma] \end{aligned}$$

whence

$$\text{best } c \left(= -\frac{\text{Cov}}{\text{Var}[\gamma]} \right)$$

Usually

unknown.