

⊗ when the dynamics are unknown, but have data:
 first look at value prediction / Policy evaluation

$$V^\pi(s) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r(s_t, a_t) \mid s_0 = s \right]$$

$$a_t \sim \pi(\cdot \mid s_t)$$

$$s_{t+1} \sim P(\cdot \mid s_t, a_t)$$

$(s_t, a_t, r_t, s_{t+1})_{t=0}^T \Rightarrow Q: \text{how to estimate } V^\pi(s) ?$

$$1) \quad V^\pi(s) = r^\pi(s) + \gamma \sum_{s' \in \mathcal{S}} P^\pi(s' \mid s) V^\pi(s')$$

If $v \in \mathbb{R}^{|\mathcal{S}| \times 1}$, then

$$V^\pi = r + \gamma P V^\pi \Rightarrow r + (\gamma P - I) V^\pi = 0$$

$$\frac{d}{dt} V(t) = \underbrace{r + (\gamma P - I) V(t)}_A \xrightarrow{t \rightarrow \infty} V(t) \rightarrow V^\pi$$

$\text{eig}(A) < 0$

(think: why $\frac{d}{dt} V(t) = (I - \gamma P) V(t) - r \xrightarrow{t \rightarrow \infty} V^\pi$).

$$2). \quad V_{n+1} = V_n + \alpha_n (r + (\gamma P - I) V_n)$$

If the current state is x_n

$$V_{n+1}(x_n) = V_n(x_n) + \alpha_n \underbrace{(r(x_n) + \gamma V(x_{n+1})) - V(x_n))}_{TD}$$

Then: $\sum_n \alpha_n = \infty \quad \sum_n \alpha_n^2 < \infty \Rightarrow V_n \rightarrow V^\pi$

TD is a stochastic Approximation method to approximate the steady state of the ODE $\dot{x} = h(x)$.

$$x_{n+1} = x_n + \varepsilon_n (h(x_n) + \frac{M_{n+1}}{\varepsilon_n})$$

$$\mathbb{E}[M_{n+1}] = 0$$

$$\Leftrightarrow x_{n+1} = x_n + \varepsilon_n \hat{h}(x_n), \text{ where } \mathbb{E}[\hat{h}(x_n) | x_n] = h(x_n)$$

with assumption on $h, \varepsilon_n, M_n, x \rightarrow x^*$
 \hookrightarrow stable equilibrium.

Robbins-Monro 1951

Assumption on h, ε_n, M_n :

① $h \rightarrow$ Lipschitz: $\exists L$ s.t. $\|h(x) - h(y)\| \leq L \|x - y\|$ for x, y

② $\sum_{n \geq 0} \varepsilon_n = \infty, \sum_{n \geq 0} \varepsilon_n^2 < \infty$

③ $\mathbb{E}[M_{n+1} | \mathcal{F}_n] = 0$ & $\mathbb{E}[\|M_{n+1}\|^2 | \mathcal{F}_n] \leq k(1 + \|x_n\|^2)$

($\Leftrightarrow \mathbb{E}[V(S_{t+1})^2 | S_t] - (\mathbb{E}[V(S_{t+1}) | S_t])^2 \leq \|V\|^2$)

④ $\sup_{n \geq 0} \|x_n\| < \infty$ (\Leftarrow By "Borkar-Meyn, 2000")

⑤ $\exists f(x)$ s.t. $\lim_{\|x\| \rightarrow \infty} f(x) = +\infty$ & $f(x) \geq 0, f(x) \in C^1 \leftarrow$ Lyapunov fun.

s.t. $\langle \nabla f(x), h(x) \rangle \leq 0$ w. strict neg if $f(x) \neq 0$

x^* is a stable equl

For linear ODE $\text{Re}(\lambda) < 0$ is enough (see the proof in the next page).

A linear dynamical system is asymptotically stable

if for \forall initial condition $x(0)$, $x(t)$ converges to the origin 0 as $t \rightarrow \infty$

① $\dot{x}(t) = Ax(t)$ is asymp stable $\Leftrightarrow |e^{j\omega(A)}| < 1$

1) If $\exists Q$ & Λ s.t. $A = Q \Lambda Q^{-1}$
 (c_1, \dots, c_n) $\text{diag}(\lambda_1, \dots, \lambda_n)$.

then $x(t) = Q \Lambda^n Q^{-1} x(0)$

$$= (c_1, \dots, c_n) \begin{matrix} \lambda_1^t y_1(0) \\ \lambda_2^t y_2(0) \\ \vdots \\ \lambda_n^t y_n(0) \end{matrix}$$

$$\|x(t)\| = \left\| \sum_{i=1}^n c_i \lambda_i^t y_i(0) \right\| \leq \sum_{i=1}^n |c_i| |\lambda_i|^t |y_i(0)| \rightarrow 0$$

2) Jordan Normal form Thm:

\forall matrix A can be conjugated to $\Lambda + N$ with $N^n = 0$

$$(B+N)^t = \begin{matrix} B^t & C_1^t B^{t-1} N & \dots & C_{n-1}^t B^{t-n+1} N^{n-1} \\ \downarrow & \downarrow & & \downarrow \\ 0 & 0 & & 0 \end{matrix}$$

② $\dot{x}(t) = Ax(t)$ is asymp stable $\Leftrightarrow \text{Re}(e^{j\omega(A)}) < 0$

$$x(t) = e^{At} x(0) = U^t x(0) \quad \text{with } U = e^A$$

$$|e^{j\omega(e^A)}| < 1 \Leftrightarrow |e^{e^{j\omega(A)}}| < 1 \Leftrightarrow \text{Re}(e^{j\omega(A)}) < 0$$

"The ODE method for conv of SA & RL"

"Borkar-Meyn Thm (2000)"

$$\text{Lip} + \text{stepsize} + \text{Var} + \left(\frac{h(r, x)}{r} \right)_{r \rightarrow \infty} = h_{\infty}(r, x) \quad \text{for } \forall x$$

$\dagger \theta$ is a globally exponentially stable equl of



$$\dot{\theta} = f_{\infty}(\theta)$$

$$\Rightarrow x_n \rightarrow x^*$$

For linear $f_{\infty}(\theta) = A\theta$, $\text{Re}(\text{eig}(A)) < 0$

(*) The convergence rate for standard linear SA method: $O\left(\frac{1}{\sqrt{t}}\right)$