

# Function approximation + Planning

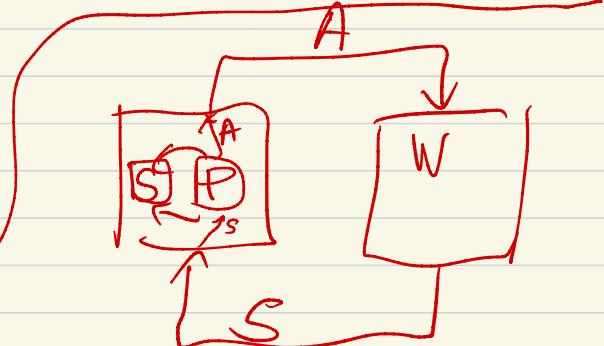
Review

Tabular MDPs

S states  
A actions

$$\rightarrow \mathcal{S}^2 (S^2 A)$$

S big!



Local Planning

$$O((H^5 A)^{\frac{H}{\delta}})$$

$$H \approx \frac{1}{1-\delta}$$

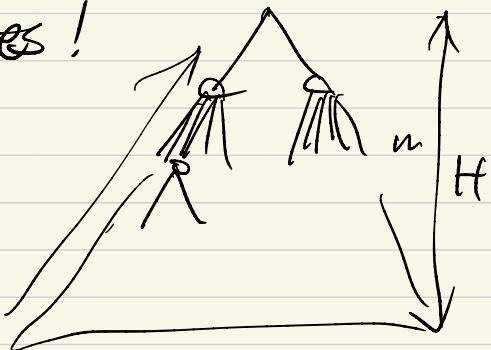
running indep. of # states!

$$\rightarrow \mathcal{S}^2 (A^{\frac{H}{\delta}})$$

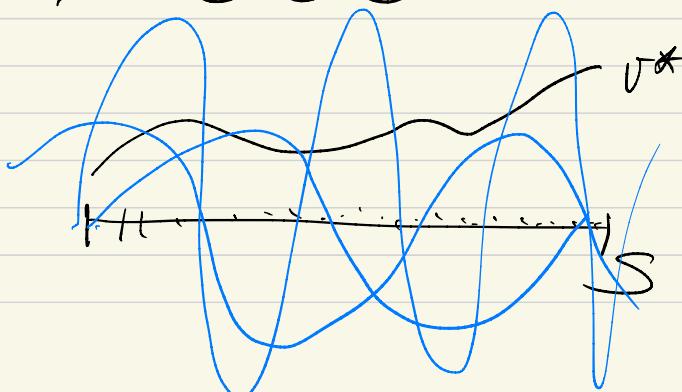
$$\tau = 0.99$$

$$1 - \delta = \frac{1}{\frac{1}{100}} = 100$$

$$A = 2, \boxed{2^{100}}$$



Value function approximation



$$V^*(S) = \sum_{i=1}^d \theta_i \varphi_i(S)$$

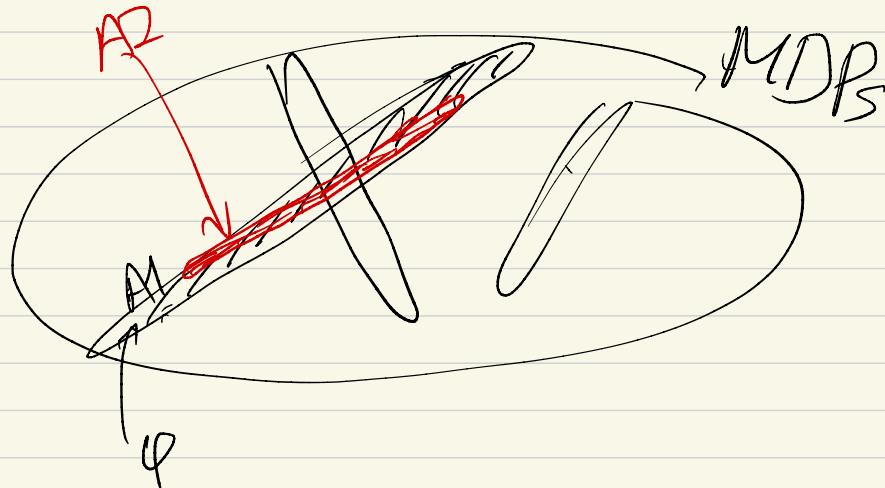
One-step lookahead

$$\hat{\pi}(s) = \arg \max \gamma a(s) + \tau \underbrace{\langle P_a(s), \hat{v} \rangle}_{\text{sampling}}$$

$$\frac{1}{1-\gamma} \approx \| \hat{v} - v^* \|_\infty \leq \underbrace{\varepsilon}_{\tau = 0, 99} \Rightarrow \frac{1}{1-\gamma} \{ v^* \geq v^* - \left[ \frac{2\tau\varepsilon}{1-\gamma} \right] \}$$

$$\theta \in \mathbb{R}^d$$

$$d \ll s$$



$$n=|S| \quad \left\{ \begin{array}{l} v^* = \begin{bmatrix} \varphi^T(s_1)\Theta_* \\ \vdots \\ \varphi^T(s_n)\Theta_* \end{bmatrix} = \begin{bmatrix} -\varphi^T(s_1) \\ \vdots \\ -\varphi^T(s_n) \end{bmatrix} \Theta_* = \Phi \Theta_* \\ \varphi(s) = [\varphi_1(s), \dots, \varphi_d(s)]^T \end{array} \right.$$

$$\begin{aligned} v^* &\in \text{span}(\Phi) = \\ &= \{ \Phi \theta \mid \theta \in \mathbb{R}^d \} \\ &\subseteq \mathbb{R}^n \end{aligned}$$

$$n \gg d$$

$$\underbrace{\mathcal{F}}_{\mathbb{R}^n}$$

~~(A1)~~  $v^* \in \text{span}(\Phi)$  :  $v^*$  - realizable

(A2)  $\nexists \pi: v^\pi \in \text{span}(\Phi)$

Note: (A2)  $\Rightarrow$  (A1); FT  $\exists \pi: v^\pi = v^*$ .

$\mathcal{P}(B1)$

$q^* \in \text{span}(\tilde{\Phi})$

$\Phi := \tilde{\Phi}$

$(B2) \quad \forall \pi: q^\pi \in \text{span}(\tilde{\Phi})$

$(B2) \Rightarrow (B1)$

$\mathcal{F} = \text{span}(\Phi)$

Today: Assume  $(B2)$  holds.

Simulator access.

3 local planners s.t.

$\text{poly}(A, d, H, \frac{1}{\delta})$  queries/compute

$\rightsquigarrow \pi: V^\pi \geq V^* - \delta I$ .

Value Iteration:  $Q_{k+1} = D$ ,  $V_{k+1} = T Q_k$

Policy Iteration:



$(B2)$

$\pi_0$  arbitrary

$\pi_{k+1}$  greedy wrt.  $V^{\pi_k}$

$\pi_{k+1}$  greedy wrt.  $Q^{\pi_k}$

$$Q^{\pi_k} = r + \gamma P V^{\pi_k}$$

$r(s,a)$

$$(Pv)(s,a) = \langle P_a(s), v \rangle$$

$$\pi_{k+1}(s) = \arg \max_a [Q^{\pi_k}(s,a)] = \arg \max_a \langle P_a(s), v \rangle$$

$$Q^{\pi_k} = \tilde{\Phi} \Theta_k, \quad \Theta_k = ?$$

$$T \mathcal{F} \subseteq \mathcal{F}$$

closedness/  
invariance

$Q \in \mathcal{F}$

$Q_{k+1} \in \mathcal{F}$

$\forall q \in \mathcal{F} \Rightarrow$

$Tq \in \mathcal{F}$

# How to calculate $\Phi\theta$ ?

$\pi: S \rightarrow \mathcal{A}$ , (B2), simulator

point-evaluate  $\pi$

Need:  $\theta \in \mathbb{R}^d$

$$\Phi\theta = q^\pi$$

Policy evaluation

$$q^\pi = r + \gamma \underbrace{P_\pi q^\pi}_{} \quad P_\pi: \mathbb{R}^{SA} \rightarrow \mathbb{R}^{SA}$$

$$q^\pi = \Phi\theta$$

$$q^\pi = r + \gamma P_\pi q^\pi$$

$$q^\pi \in \mathbb{R}^{SA}$$

SA unknowns

$$\theta \in \mathbb{R}^d$$

d unknowns

SA+d

$$\Phi\theta = r + \gamma P_\pi \Phi\theta$$

unknowns: d

$$r \in \mathbb{R}^{SA}$$

OOPS!

we have SA equations!

$$\begin{matrix} \uparrow & \\ SA & \left\{ \begin{bmatrix} I & P_\pi \end{bmatrix} \Phi \right\} \begin{bmatrix} \Theta \\ r \end{bmatrix} = \begin{bmatrix} r \\ \vdots \\ r \end{bmatrix}_{SA} \end{matrix}$$

SVD

LU decomposition

\* Subsample rows

\* Least-squares

What if (B2) does not even hold??

Robust

$\text{poly}(d, A, H, \frac{1}{\delta}, \dots)$

[indep. of  $S$ ]

Weighted least-squares.

$$\Phi \theta = q^\pi$$

$$(s, a) \in C \\ \subseteq S \times A$$

$$q^\pi \approx \sum_{t=0}^k \gamma^t r_t$$

$$\gamma^k \geq H \gamma^k \delta \Rightarrow \gamma \leq \delta$$

$$q^\pi(s, a) = \Phi(s, a)^\top \theta$$

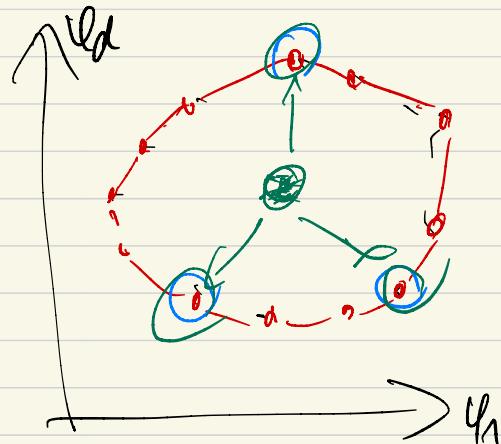
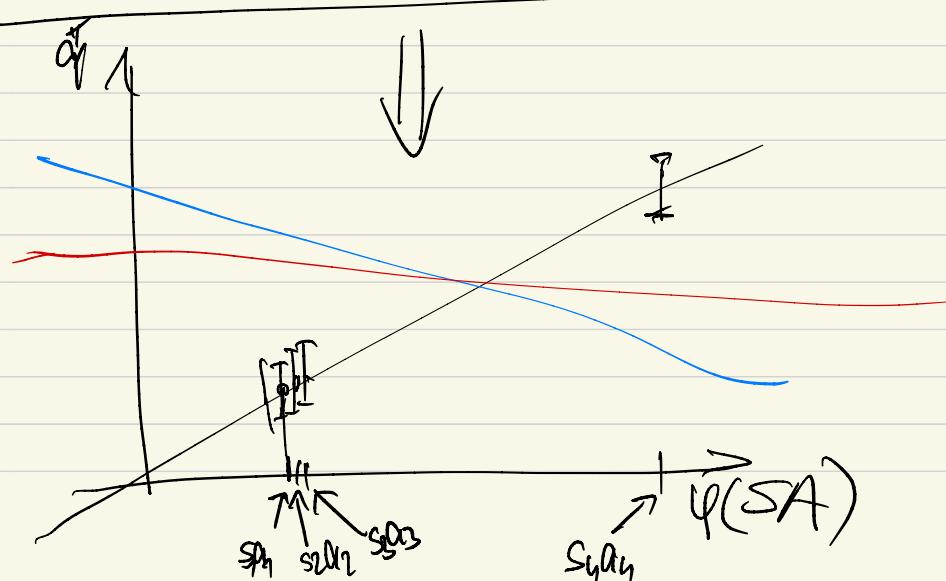
"Least-squares" / penan. wt. problem.

$$R(s, a) \approx \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^k \gamma^t r_{A_t^{(i)}}(S_t^{(i)}) = \hat{R}(s, a)$$

Simulations

$$\begin{cases} A_t^{(i)} = \pi(S_t^{(i)}), t \geq 1 \\ S_{t+1}^{(i)} \sim P_{A_t^{(i)}}(S_t^{(i)}) \\ S_0^{(i)} = s, A_0^{(i)} = a \end{cases}$$

$$\mathbb{E}[\hat{R}(s, a)] = (\Gamma^\pi \theta)(s, a)$$



Theorem (Kiefer - Wolfowitz):

$$\varphi: S \times A \rightarrow \mathbb{R}^d$$

$$\sum_{(s,a)} g(s,a) = 1$$

$$C \subseteq S \times A, g: C \rightarrow [0,1]$$

$$|C| \leq \boxed{d(d+1)/2}$$

s.t.

$$\max_{(s,a)} \boxed{\|\varphi(s,a)\|_{G_g^{-1}}} \leq \boxed{\sqrt{d}}$$

$$G_g = \sum_{\substack{(s,a) \in C}} g(s,a) \varphi(s,a) \varphi(s,a)^T$$

$$2\sqrt{d}$$

$$\|x\|_P^2 = x^T P x$$

$$|C| \approx 2d \log 6dd + 16$$

Optimal Experimental Design

(Todd, 2016)

How to find a small  $C$ ?