

March 9

# Policy Search $\rightarrow$ Policy Gradients

What to optimize?  $\rightarrow$  MDP

Value functions!

Policies

$$M = (S, A, P, r)$$

$$0 \leq \gamma < 1$$

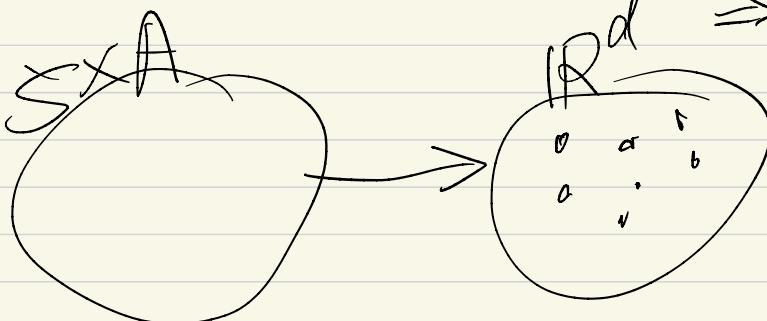
$|S|$  too big!

$$\pi \in \Pi$$

$\leftarrow$  set of memoryless policies of  $M$

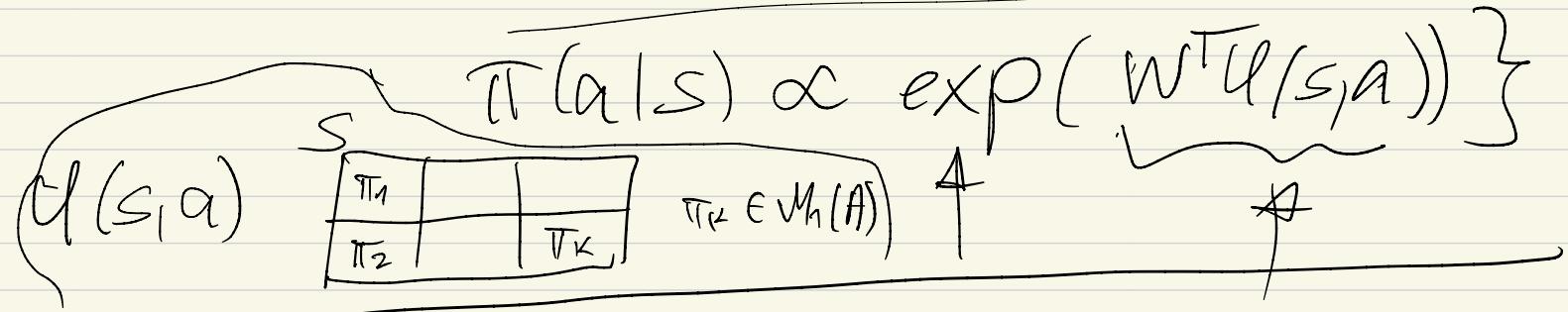
$$q: S \times A \rightarrow \mathbb{R}^d$$

$$\Pi_q = \left\{ \pi \in \Pi \mid \forall (s, a), (s', a') \in S \times A \right.$$
  
$$\left. q(s, a) = q(s', a') \right\}$$



$$\mathbb{R}^d \Rightarrow \pi(a | s) = \pi(a' | s')$$

$$\mathcal{P}_{\text{Boltzmann-Lin}} = \overline{\{ \pi \in \text{ML} \mid \exists w \in \mathbb{R}^d }$$



$$\mu \in \mathcal{M}_1(S) \quad \mu = \text{row-vector}$$

$$J(\pi) = \mu \cdot v^\pi = \sum_{s \in S} \mu(s) v^\pi(s)$$

$\underset{\pi \in \mathcal{P}}{\operatorname{argmax}}$   $J(\pi)$

"policy search"

Boltzmann-linear

State-aggregation

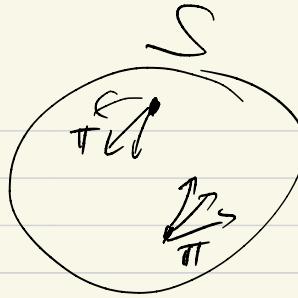
No! Vlassis Littman Barber (POMDR)

How come?

BPS = Blind-Policy-Search

$$\varphi(s, a)_{a'} = \mathbb{I}(a = a')$$

$$\mathcal{T} = \underbrace{\mathcal{M}_1([A])}_{\pi \in \mathcal{M}_1([A])}$$



$$S = A = [n]$$

$$P_a(s, a) = 1 \quad / \quad P_a(s, s') = 0, \quad s' \neq a$$

$$r_a(s) = ?$$

$$r_{\pi}(s) = \sum_a \pi(a) r_a(s)$$

$$M = \frac{1}{n} \mathbf{1}^T$$

$$J(\pi) = \mu \pi^T = \mu \left( \sum_{t=0}^{\infty} \pi^t P_{\pi}^t \right) r_{\pi}$$

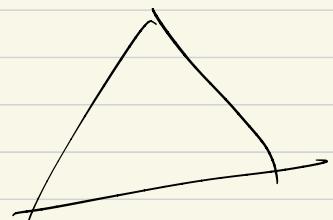
$$P_{\pi}(s, s') = \sum_{a \in A} \underbrace{\pi(a|s)}_{\pi(a)} \underbrace{P_a(s, s')}_{\mathbb{P}(s'|a)} = \pi(s')$$

$$P_{\pi} = \mathbf{1} \pi^T, \quad t \geq 1$$

$$P_{\pi}^2 = \mathbf{1} \underbrace{\pi^T}_{1} \mathbf{1} \pi^T = \mathbf{1} \pi^T \quad \left| \quad R = (r_a(s))_{s,a} \right.$$

$$P_{\pi}^t = \mathbf{1} \pi^T$$

$$J(\pi) = \underbrace{\mu \left( \mathbf{I} + \frac{\sigma}{1-\delta} \mathbf{1} \pi^T \right) R}_{x} \pi$$



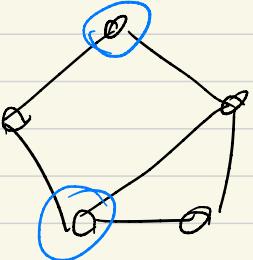
$$= \boxed{nR\pi + \frac{\gamma}{1-\gamma} \cdot \underline{\pi^T R \pi}} \rightarrow \max_{\pi \in M_1(A)}$$

## Graph Theory

### MAX-INDSET

$\Delta(G)$

$G = (V, E)$



indep set = neighbor-free set.

NP-hard

$$\Delta(G) = \max \{ |V'| \mid V' \subseteq V, \text{indep. in } G \}$$

cubic (3-regular)

Theorem (Motzkin-Strauss '65) :

$$|V|=n \quad \frac{1}{\Delta(G)} = \min_{y \in M_1(n)} \underline{\underline{y^T (G+I) y}}$$

$\{0,1\}^{n \times n} \ni G$  vertex adj. matrix of  $G = (V, E)$

$$G_{i,j} = 1 \iff (i, j) \in E, i \neq j$$

Pick  $G$  3-regular

$$R := -(I + G)$$

$$J(\pi) = -\mu(I + G)\pi - \frac{\gamma}{1-\gamma} \pi^T (I + G)\pi$$

$$= \frac{1}{n} \underline{1^T (\mathbb{I} + G) \pi} - \frac{\gamma}{1-\gamma} \pi^T (\mathbb{I} + G) \pi$$

$$= \frac{1}{n} - \frac{\gamma}{1-\gamma}$$

$$(1^T G)_j = \sum_i G_{ij} = 3, \quad 1^T G = 3 \quad 1^T$$

$$\max_{\pi \in \Pi} J(\pi) = \frac{1}{n} - \frac{\gamma}{1-\gamma} \frac{1}{\alpha(G)} \geq -\frac{1}{n} - \frac{\gamma}{1-\gamma} \frac{1}{M}$$

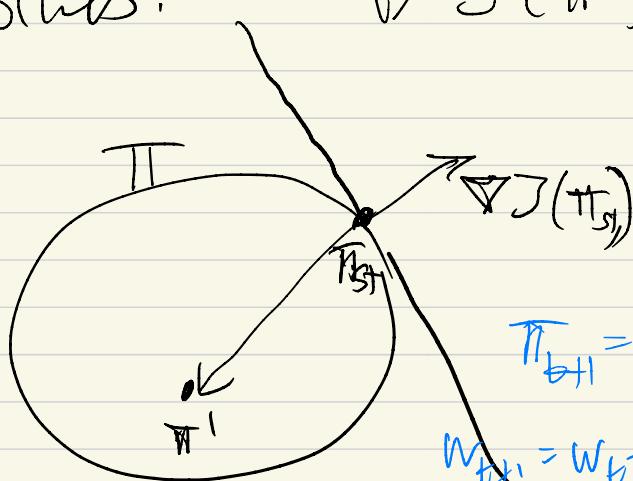
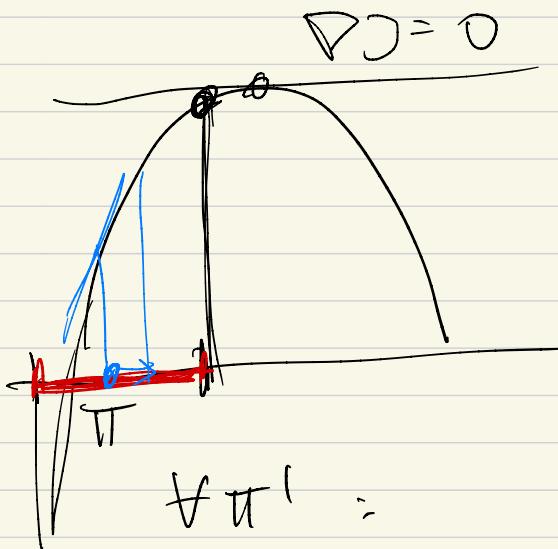
$$\Leftrightarrow \alpha(0) \geq m$$

Conclusion: Policy Search can be hard!

Relax the problem!

Stationary points!

$$\nabla J(\pi) = 0$$



$$\forall \pi': \quad \langle \nabla J(\pi_{st}), \pi' - \pi_{st} \rangle \leq 0$$

# 1. Computation (nice & benign)

## 2. Approx. guarantees

$\pi_{st}$  stat. point

$$\Rightarrow \mu v^T \pi_{st} \geq \mu v^T \pi^* - \underbrace{?}_{\epsilon(M, T)} \\ \epsilon(M, \varphi)$$

gradient ascent

1.1  $\nabla_w J(\pi_w) = ?$  simulator?

1.2 Which algo?  $\xrightarrow{\text{Agnostic to } \pi}$   $\xrightarrow{\text{Nope!}} ?$

choose vanilla PG method

Natural-PG

$$\nabla_w J(\pi_w) = \nabla_\pi J(\pi_w) \nabla_w \pi_w$$

