

March 11

# Gradient Computation

PG - Theorem

Bhandari - Russo

Silver

Sutton / Safrader / Kakade / ..

Deterministic policies

$$\pi_\theta : S \rightarrow \mathcal{A}$$

$$A \subseteq \mathbb{R}^F$$

Finite-action MDPs

$$\tilde{\mathcal{A}} = \mathcal{M}_1(A) = \{p \in [0, 1]^A \mid \sum_{a=1}^A p_a = 1\}$$

$$\underline{M} = (S, P, r)$$

$$(P, r)$$

$$P = (P_a(s))_{a \in A}$$

$$\begin{cases} \tilde{\pi} : S \rightarrow \mathcal{M}_1(\mathcal{A}) \\ \pi : S \rightarrow \tilde{\mathcal{A}} \end{cases}$$

$$\tilde{P} = (\tilde{P}_{\tilde{a}})_{\tilde{a} \in \tilde{\mathcal{A}}}$$

$$\tilde{P}_{\tilde{a}}(s, s') = \sum_{i=1}^A \tilde{a}_i P_i(s, s')$$

$$\tilde{r} = (r_{\tilde{a}})_{\tilde{a} \in \tilde{\mathcal{A}}}$$

$$r_{\tilde{a}}(s) = \sum_{i=1}^A \tilde{a}_i r_i(s)$$

$$\underline{\tilde{M}} = (S, \tilde{\mathcal{A}}, \tilde{P}, \tilde{r})$$

$\pi$  SML of  $M$

$\tilde{\pi}$  DML of  $\tilde{M}$

$$\tilde{\pi}(s) = \pi(s)$$

$$v^\pi = \tilde{v}^\pi$$

$$v^* = \tilde{v}^*$$

$$q^\pi = \tilde{q}^\pi$$

det. policies of  $\tilde{M}$  suffice for  $M$ .

$$A \subseteq \mathbb{R}^P, \quad \pi_\theta : S \rightarrow A, \quad \theta \in \mathbb{R}^d, \quad x \in \mathbb{R}^d$$

Theorem:  $f_{\pi} : x \mapsto \sum_{\pi} \tilde{\mu}_{\pi} M_{\pi_x} q^{\pi}$

$$\boxed{(M_{\pi} q)(s) = q(s, \pi(s))}$$

$$g_{\pi} : x \mapsto \tilde{\mu}_{\pi_x} v^{\pi}$$

$$\text{Fix } \theta_0 \in \mathbb{R}^d, \quad \mathcal{J}(\theta) = \mu v^{\pi_{\theta}}$$

Assume  $f_{\pi_{\theta_0}}, g_{\pi_{\theta_0}} \in C^1 @ \theta_0$

$$\begin{aligned} \Rightarrow \nabla \mathcal{J}(\theta_0) &= \frac{d}{dx} \underbrace{\tilde{\mu}_{\pi_{\theta_0}} M_{\pi_x} q^{\pi_{\theta_0}}}_{x=\theta_0} \\ &= \frac{d}{dx} \sum_s \tilde{\mu}_{\pi_{\theta_0}}(s) q^{\pi_{\theta_0}}(s, \pi_x(s)) \Big|_{x=\theta_0} \end{aligned}$$

$$f(u, v) \quad \frac{d}{dx} f(x, x) = \underbrace{\frac{\partial}{\partial u} f(x, x)}_{\uparrow} + \underbrace{\frac{\partial}{\partial v} f(x, x)}_{\uparrow}$$

Proof:  $v^{\pi'} - v^{\pi} = (\mathbb{I} - \gamma P_{\pi'})^{-1} r_{\pi'} - v^{\pi}$

$$\begin{aligned} &= (\mathbb{I} - \gamma P_{\pi'})^{-1} [r_{\pi'} - (\mathbb{I} - \gamma P_{\pi'}) v^{\pi}] \\ &= (\mathbb{I} - \gamma P_{\pi'})^{-1} [T_{\pi'} v^{\pi} - v^{\pi}] \end{aligned}$$

$$T_{\pi'} r^{\pi} = M_{\pi'} (\underbrace{r + \gamma P_{\pi'} v^{\pi}}_{q^{\pi}})$$

$$M_{\pi'} q^{\pi}$$

$$M(V^{\pi_x} - \underline{V^{\pi_0}}) = \underbrace{\tilde{P}_\mu^{\pi_x} [M_{\pi_x} q^{\pi_0} - V^{\pi_0}]}_{\text{/ } \frac{d}{dx}}$$

$$f(u, v) = \tilde{P}_\mu^{\pi_u} M_{\pi_v} q^{\pi_0}$$

$$\begin{aligned} \frac{d}{dx} f(x, x) &= \left. \frac{d}{du} \tilde{P}_\mu^{\pi_u} M_{\pi_x} q^{\pi_0} \right|_{u=x} \\ &\quad + \left. \frac{d}{dv} \tilde{P}_\mu^{\pi_x} M_{\pi_v} q^{\pi_0} \right|_{v=x} \end{aligned}$$

$$\frac{d}{dx} J(x) = \left. \frac{d}{du} \tilde{P}_\mu^{\pi_u} M_{\pi_x} q^{\pi_0} \right|_{\substack{u=x \\ v=\theta}} + \left. \frac{d}{dv} \tilde{P}_\mu^{\pi_x} M_{\pi_v} q^{\pi_0} \right|_{v=x} - \left. \frac{d}{dx} \tilde{P}_\mu^{\pi_x} V^{\pi_0} \right|_{x=\theta}$$

$$x = \theta = \theta_0$$

$$\nabla J(\theta_0) = \boxed{\left. \frac{d}{dv} \tilde{P}_\mu^{\pi_{\theta_0}} M_{\pi_v} q^{\pi_{\theta_0}} \right|_{v=\theta_0}} // \text{Qm. c.d.}$$

$$\pi_\theta : S \rightarrow M_1 / \mathcal{A}; \quad \pi_\theta(a | s)$$

$$\begin{aligned} \frac{d}{dv} \tilde{P}_\mu^{\pi_\theta} M_{\pi_v} q^{\pi_\theta} \Big|_{v=\theta} &= \frac{d}{dv} \sum_s \tilde{P}_\mu^{\pi_\theta}(s) \sum_a \pi_v(a | s) q^{\pi_\theta}(s, a) \\ &\quad \xrightarrow{\text{b}(s)} \end{aligned}$$

$$= \sum_s \tilde{P}_\mu^{\pi_\theta}(s) \sum_a \frac{d}{dv} \pi_v(a | s) \Big|_{v=\theta} \underbrace{q^{\pi_\theta}(s, a)}$$

$$\sum a(s) = 1$$

$$\frac{d}{dw} \log \pi_v(a|s) = \frac{\frac{d}{dw} \pi_v(a|s)}{\pi_v(a|s)}$$

$$\nabla J(\theta) = \sum_s \tilde{p}_\mu^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) q^{\pi_\theta}(s, a) \frac{\partial}{\partial w} \log \pi_v(a|s)$$

$\uparrow$        $\uparrow$        $\uparrow$

$\pi_v(a|s) > 0$

$v = \theta$

Simulation to get  $\hat{\nabla} J(\theta)$

$$(1-\gamma) \tilde{p}_\mu^{\pi_\theta}(s) \pi(a|s) = (1-\gamma) \underbrace{p_\mu^{\pi}(s, a)}$$

$$\sum_{s_i \in \mathcal{S}} = 1$$

(1)  $(s_i, a_i) \sim (1-\gamma) p_\mu^{\pi_\theta}(\quad), \quad i=1, \dots, n$

(2) Rollout from  $s_i, a_i$  with  $\pi_\theta$

) for  $H_\theta, \epsilon(1-\gamma)$  for  
 $\epsilon$  error,  $\hat{Q}_{i,j}$   $j=1, \dots, m$   
 $i$  indep.  
 $N \sim \text{Geo}(\gamma)$

$$\nabla J(\theta) = \frac{1}{(1-\sigma)n^m} \sum_{i=1}^n \sum_{j=1}^m Q_{ij} \frac{\partial \log \pi_\theta(A_i|S_j)}{\partial \theta}$$

Vanilla PG : too bad. / Frank-Wolfe

NPG : Natural Policy Gradient

$$F(\theta) = \nu_n^{\pi_\theta} \left( \nabla_\theta \log \pi_\theta(\cdot | \cdot) \quad \nabla_\theta \log \pi_\theta(\cdot | \cdot)^T \right) \in \mathbb{R}^{d \times d}$$

$$\text{NPG: } \theta_{t+1} = \theta_t + \gamma_t \underbrace{F(\theta_t)^T \nabla_\theta J(\theta_t)}$$

$$\text{Thm: } (1-\sigma) F(\theta)^T \nabla_\theta J(\theta) =$$

$$= \underset{\substack{w \in \mathbb{R}^d \\ \min \cdot \text{norm}(w, \|\cdot\|_2)}}{\arg \min} \nu_n^{\pi_\theta} \left( w^T \underbrace{\nabla_\theta \log \pi_\theta(a|s)}_{\in \mathbb{R}^d} - a^{\pi_\theta} \right)^2$$

$$a^{\pi_\theta} = q^{\pi_\theta} - v^{\pi_\theta}$$

Proof: ✓

$$\pi_\theta(a|s) \propto \exp(-\theta^\top \varphi(s, a))$$

$$\nabla \log \pi_\theta(a|s) = \varphi(s, a) - \underbrace{\sum_{a'} \pi_\theta(a'|s) \varphi(s, a)}_{\psi_\theta(s, a)}$$

$$NPG: \quad \Delta \theta_t = \gamma w_t \quad w_t = \underset{w}{\operatorname{argmin}} \sum_m^M [w^\top \varphi - a^{\pi_\theta}]^2$$

$\downarrow$   
Tabular;  $\varphi \neq \text{tabular}$

$$\left\| \frac{d\varphi}{da} \right\|_\infty = C$$

$$\boxed{Q-NPG: \quad \Delta \theta_t = \gamma w_t \quad w_t = \underset{w}{\operatorname{argmin}} \sum_m^M [w^\top \varphi - q^{\pi_\theta}]^2}$$

Polifex

$$\Delta \theta_t = \gamma w_t \quad w_t = \underset{w}{\operatorname{argmin}} \sum_m^M [w^\top \varphi - q^{\pi_\theta}]^2$$

Mirror-descent

$$\frac{\nabla d_E}{1-\gamma}$$

$$\frac{C_E}{(1-\gamma)^{1/2}}$$

Questions: ① What parametrization?

Escort map?

② Adaptive directions/stepsizes

$$\begin{aligned} s &\in \mathbb{R}^A \\ |s_i|^{1/2} \\ \sum_j |s_j|^{1/2} \end{aligned}$$