

Lecture 26

CS 59000 - RL Theory

Policy Gradient

Agenda

- Performance Difference Lemma
- Gradient Domination
- Global optimality ☺

For discounted setting

$$\nabla_{\theta} J(\pi_{\theta}) = \int_{\mathcal{X}} \int_{\mathcal{X}} \int_{\mathcal{A}} \underbrace{\mu_{\pi_0}^{\delta}(x', u)}_{\text{occupancy kernel}} \nabla_{\theta} \pi(a; u) Q_{\pi_0}(u, a) da) dx' dx$$

For undiscounted.

$$\nabla_{\theta} J(\pi_{\theta}) = \int_{\mathcal{X}} \mu_{\pi_{\theta}}(u) \left(\int_{\mathcal{A}} \nabla_{\theta} \pi(a; u) Q_{\pi_{\theta}}(u, a) da \right) du$$

⇒ ways to estimate: \int Monte Carlo sampling \rightarrow Alexandrov 1965
Bayesian Quadrature

→ Bayesian Quadrature Policy Gradient
AKella, 2020

Parametrization ; $\theta \in \Theta$; π_θ

Consider the finite state action MDP.

- Direct parametrization

$$\rightarrow \pi_\theta(a; n) = \theta_{an} \geq 0 \text{ and } \sum_a \theta_{an} = 1 \quad \forall n \in \mathcal{X}$$

- Softmax parametrization

$$\pi_\theta(a; n) = \frac{\exp(\theta_{an})}{\sum_{a'} \exp(\theta_{a'n})}$$

(side note $\nabla_{\theta} \eta(\pi) \rightarrow \theta \neq \nabla \eta(\pi) \Rightarrow$ need to project onto simplex)
For direct parametrization

General Function Classes

$$\pi_\theta(a; n) = f_\theta(n, a) ; \sum_a f_\theta(n, a) = 1$$

Consider the discounted finite state action MDP.

Define: Advantage function $A_{\pi}(x, a); \forall x, a \in \mathcal{X} \times \mathcal{A}$

$$A_{\pi}(x, a) = \underbrace{Q_{\pi}(x, a)} - \underbrace{V_{\pi}(x)}$$

Lemma (Performance difference Lemma; Kakade & Langford 2002)

$$\eta(\pi) - \eta(\pi') = E_{\mu_{\pi}^{\gamma}} [E_{\pi} [A_{\pi'}(x, A) | x]]$$

proof in Final.

(Azizzadenesheli 2018)
POMDP et al

What is policy gradient for direct parametrization

$$\underbrace{\mu_{\pi}^{\gamma}(x')} = \sum_x \underbrace{\mu_{\pi}^{\gamma}(x'; x)} \underbrace{P_1(x)}$$

$$\underbrace{\nabla_{\theta} \eta(\pi_{\theta})}_{=} = (1-\gamma) \sum_{x'} \sum_a \nabla_{\theta} \pi_{\theta}(a; x) Q_{\pi}(x', a) \mu_{\pi_{\theta}}(x') da dx$$

$$\frac{\partial \eta(\pi_{\theta})}{\partial \theta_{an}} = (1-\gamma) Q_{\theta}(x, a) \mu(x) \leftarrow$$

$$\rightarrow \nabla \eta(\pi_{\theta}) = (1-\gamma) \mu Q$$

Let π^* denote an optimal policy

Lemma: For direct parametrization

$$\frac{\eta(\pi^*) - \eta(\pi)}{1-\gamma} \leq \left\| \frac{\mu_{\pi^*}^\gamma}{\mu_\pi^\gamma} \right\|_\infty \max_{\pi'} (\pi - \pi')^\top \nabla_\theta \eta(\pi_0)$$

proof:

Gradient dominance setting

$$\frac{\eta(\pi^*) - \eta(\pi)}{1-\gamma} = \sum_{n, a \in \mathcal{X} \times \mathcal{A}} \mu_{\pi^*}^\gamma(n) \pi^*(a; n) A_\pi(n, a)$$

$$\leq \sum_{n \in \mathcal{X}} \mu_{\pi^*}^\gamma(n) \max_{a'} A_\pi(n, a')$$

Implicit assumption $\left\| \frac{\mu_{\pi^*}^\gamma}{\mu_\pi^\gamma} \right\|_\infty < \infty$

$$= \sum_{n \in \mathcal{X}} \underbrace{\frac{\mu_{\pi^*}^\gamma(n)}{\mu_\pi^\gamma(n)}}_{\text{Implicit assumption}} \underbrace{\mu_\pi^\gamma(n)}_{\text{Non-negative}} \max_{a'} A_\pi(n, a')$$

$$\leq \left\| \frac{\mu_{\pi^*}^\gamma}{\mu_\pi^\gamma} \right\|_\infty \sum_{n \in \mathcal{X}} \mu_\pi^\gamma(n) \max_{a'} A_\pi(n, a')$$

$$= \left\| \frac{\mu_{\pi^*}^r}{\mu_{\pi}^r} \right\| \max_{\pi'} \sum_{(u,a) \in \mathcal{X}_{\text{ste}}} \mu_{\pi}^r(u) \pi'(a; u) A_{\pi}(u, a)$$

Remember $\sum_a \pi(a; u) A_{\pi}(u, a) = 0$

$$\Rightarrow \left\| \frac{\mu_{\pi^*}^r}{\mu_{\pi}^r} \right\| \max_{\pi'} \sum_{(u,a) \in \mathcal{X}_{\text{ste}}} \mu_{\pi}^r(u) \underbrace{(\pi'(a; u) - \underbrace{\pi(a; u)}_{Q_{\pi}(u, a)})}_{Q_{\pi}(u, a) - V(u)} A_{\pi}(u, a)$$

$$\Rightarrow \left\| \frac{\mu_{\pi^*}^r}{\mu_{\pi}^r} \right\| \max_{\pi'} \sum_{(u,a) \in \mathcal{X}_{\text{ste}}} \mu_{\pi}^r(u) \underbrace{(\pi'(a; u) - \pi(a; u))}_{Q_{\pi}(u, a) - V(u)} \underbrace{Q_{\pi}(u, a)}_{Q_{\pi}(u, a)}$$

$$= \left\| \frac{\mu_{\pi^*}^r}{\mu_{\pi}^r} \right\| \max_{\pi'} \underbrace{\nabla_{\pi} \eta(\pi)^T}_{\nabla_{\pi} \eta(\pi)^T} \underbrace{(\pi' - \pi)}_{(\pi' - \pi)}$$

\Rightarrow Projected gradient ascent:

$$t=0 \quad ; \quad \pi_0 = \pi_{\text{init}}$$

$$\pi_{t+1} = \text{Proj}_{\Delta(A)^{|\mathcal{X}|}} \left(\pi_t + \alpha \nabla_{\pi} \eta(\pi) \Big|_{\pi=\pi_t} \right)$$

Theorem: (Projected gradient ascent);

Agarwal 2019
The projected gradient ascent algorithm on $\eta(\pi)$ with stepsize $\alpha = \frac{(1-\gamma)^3}{2\gamma|A|}$, runs for

T iterations satisfies,

$$\min_{t \leq T} \eta(\pi^*) - \eta(\pi_t) \leq \varepsilon$$

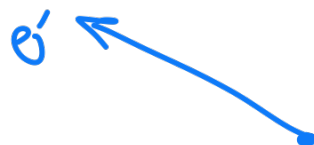
$$\text{whenever } T > \frac{64 \gamma |X| |A|}{(1-\gamma)^6 \varepsilon^2} \sup_{\pi} \left\| \frac{\mu_{\pi}^{\perp}}{\mu_{\pi}} \right\|$$

Proof Final.

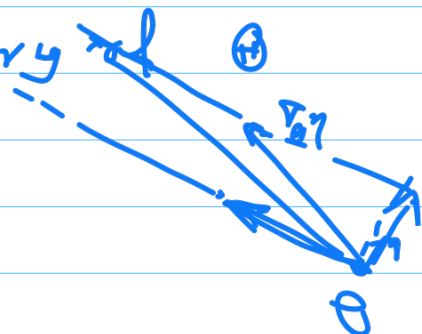
Natural policy gradient

Let's imagine our policy kernel is parametrized with θ ; π_{θ} from $(x, \omega(x))$ to $(a, \omega(a))$

remember $\eta(\pi_{\theta}) = \int_x \mu_{\pi}(x) \left(\int_a \underbrace{\pi_{\theta}(a; x)}_{\theta} \bar{r}(x, a) da \right) dx$



- what if we change θ , but Π_θ would not change!
- we want to change θ such that Π_θ changes sufficiently
- ⇒ we change the geometry



Consider Fisher information (state-dependent)

$$F(x; \theta) = E_{\Pi_\theta} \left[\nabla_\theta \log \Pi_\theta \nabla_\theta \log \Pi_\theta^T | x \right]$$

and Fisher information matrix

$$F(\theta) = E_{\Pi_\theta} \left[F(x; \theta) \right]$$

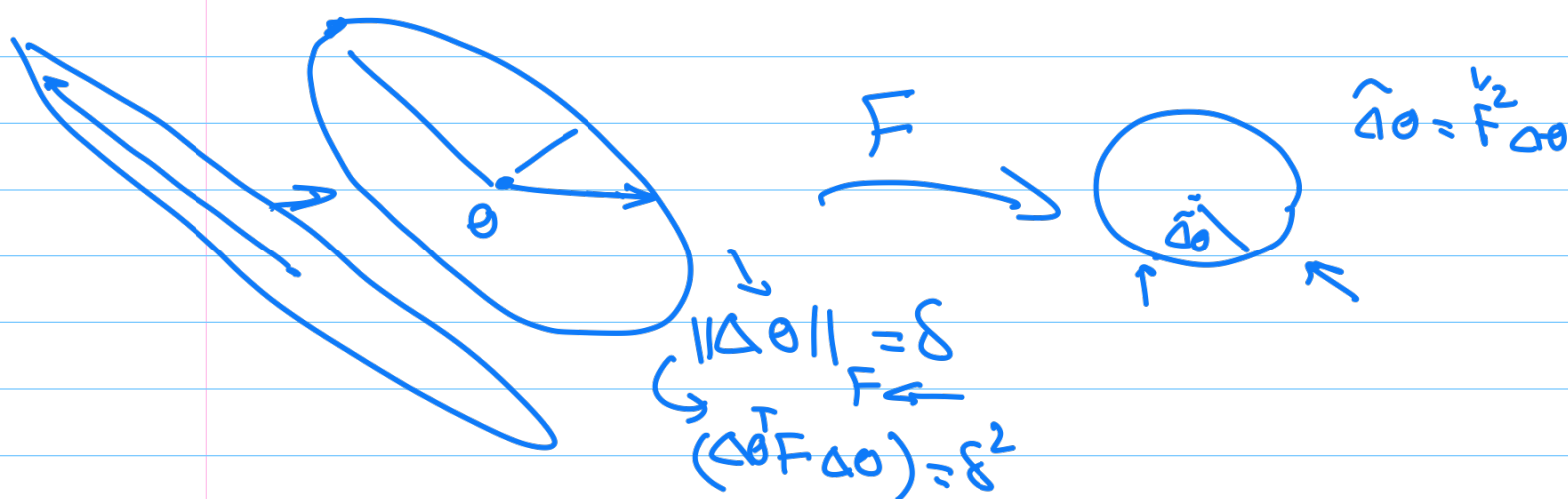
Fisher information matrix induces a Riemannian geometry such that locally the parameters of Π are metric invariant

(look Amari's book) information Geometry and application 1995 et al.)

Monday, November 23, 2020

\Rightarrow Natural policy gradient:

$$\tilde{\nabla}_{\theta} \eta(\pi_{\theta}) = \underbrace{F(\theta)^{-1}} \nabla \eta(\theta)$$



$$KL(\pi_{\theta}, \pi_{\theta + \Delta\theta}) \approx \|\Delta\theta\|_F$$

Natural gradients are essential in probabilistic optimization and the design of geometry-based optimizers.