

Lecture 27

CSS9006-RL

Policy Gradient

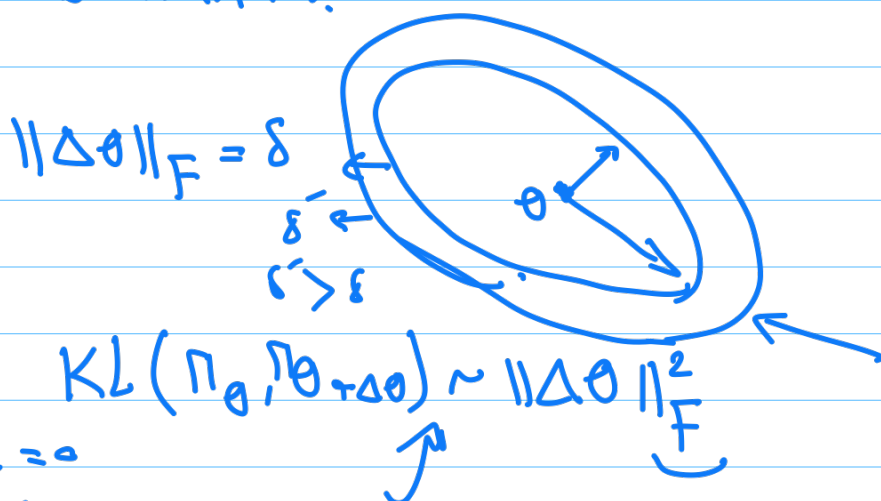
Agenda

- Natural policy gradient
- Trust region policy gradient
- Linear dynamical systems.

$$F(\chi; \theta) = E_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} \nabla_{\theta} \log \pi_{\theta}^T]$$

$$\Rightarrow F(\theta) = E_{\mu_{\pi_{\theta}}} [F(\chi; \theta)] \Rightarrow \text{Fisher information}$$

Fisher information matrix induces a Riemannian geometry such that, locally, the parameters of π_{θ} are metric invariant.



$$\Delta \theta \approx 0 \Rightarrow KL(\pi_{\theta}, \pi_{\theta + \Delta \theta}) \sim \|\Delta \theta\|_F^2$$

$$\Delta \theta = 0 \Rightarrow KL = 0$$

$$\Delta \theta \approx 0 \Rightarrow \nabla KL = 0$$

$$\Delta \theta = 0 \Rightarrow \nabla^2 KL = F$$

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

$F(\theta)$ is huge and hard to invert.

Can we compute $\tilde{\nabla} \eta(\theta) = F(\theta)^{-1} \nabla \eta(\theta)$ directly without computing $F(\theta)$ and $\nabla \eta(\theta)$ separately.

Consider a feature vector $\psi_\pi(n, a) = \nabla \log \pi_\theta(a; n)$

we use this feature vector to approximate Q .

$$\Rightarrow \omega^T \psi \sim Q \quad \text{i.e.} \quad (\omega^T \psi)(n, a) \sim Q_\pi(n, a)$$

Let's define

$$\Sigma(\omega, \pi) = \sum_x \sum_n \mu_{\pi_\theta}(n) \pi_\theta(a; n) \left((\omega^T \psi)(n, a) - Q_\pi(n, a) \right)^2$$

Theorem (Kakade 2002 "A Natural Policy gradient")

Let $\tilde{\omega}$ be a minimizer of $\Sigma(\omega, \pi)$

$$\Rightarrow \text{Then } \tilde{\omega} = \tilde{\nabla} \eta(\theta)$$

Proof: $\nabla_{\omega} \Sigma(\omega, \pi) = 0$

$$\Rightarrow \sum_{x,a} \mu_{\theta}(x) \pi_{\theta}(a|x) \psi_{\pi_{\theta}}(x,a) \psi_{\pi_{\theta}}(x,a)^T \tilde{\omega} = \sum_{x,a} \mu_{\theta}(x) \pi_{\theta}(a|x) \psi_{\pi_{\theta}}(x,a) \psi_{\pi_{\theta}}(x,a)^T \tilde{\omega}$$

Remember $\nabla_{\theta} \pi_{\theta}(a|x) = \pi_{\theta}(a|x) \nabla \log \pi_{\theta}(a|x)$

$$= \pi_{\theta}(a|x) \psi_{\pi_{\theta}}(x,a)$$

$$\Rightarrow \sum_{x,a} \mu_{\theta}(x) \pi_{\theta}(a|x) \psi_{\pi_{\theta}}(x,a) \psi_{\pi_{\theta}}(x,a)^T \tilde{\omega} = \nabla \eta(\theta)$$

$$\Rightarrow F(\theta) \tilde{\omega} = \nabla \eta(\theta) \Rightarrow \tilde{\omega} = F(\theta)^{-1} \nabla \eta(\theta)$$

$$F = \nabla_{\theta}^2 \text{KL}(\theta, \theta')|_{\theta=\theta}$$

$$\Rightarrow \eta(\pi') - \eta(\pi) \stackrel{\downarrow}{\Rightarrow} \mu_{\pi'} \left[E_{\pi} [A_{\pi}(x,a) | x] \right]$$

If we follow policy π , we can estimate A_{π}
 we have trajectories under π too

Can we optimize $\eta(\pi')$ to find an optimal policy?

If we have trajectories out of $\mu_{\pi'}$, then
 we could. we have trajectories induced by μ_{π} .

We can define and compute a surrogate function

$$\underbrace{L_{\Pi}(\Pi')} = \underbrace{E_{\Pi} \left(E_{\Pi'} [A_{\Pi}(x, A) | x] \right)} + \eta(\Pi)$$

If we keep Π' close to Π , then $L_{\Pi}(\Pi')$ is a good surrogate of $\eta(\Pi')$.

$$\text{For } \Pi' = \Pi \rightarrow L_{\Pi}(\Pi) = \eta(\Pi)$$

$$\text{and } \nabla_{\Pi'} L_{\Pi}(\Pi') \Big|_{\Pi' = \Pi} = \nabla \eta(\Pi)$$

$\Rightarrow L_{\Pi}(\Pi')$ is equal to $\eta(\Pi')$ up to first order in the vicinity of Π .

Great news:

Theorem (Schulman et al 2017) (Kakade & Langford 2002)

$$\eta(\Pi') \geq L_{\Pi}(\Pi') - \frac{4 \Sigma_{\Pi} \gamma}{(1-\gamma)^2} \max_n KL(\Pi(\cdot|n), \Pi'(\cdot|n))$$

$$\text{where } \Sigma_{\Pi} = \max_{x, a} |A_{\Pi}(x, a)|$$

[Monotonic Improvement Lemma]

Monday, November 30, 2020

optimize the $L_n(\pi) - \frac{4\varepsilon_n \gamma}{(L\gamma)^2} \max_n KL(\pi(\cdot; u), \pi'(\cdot; u))$

when $\pi' = \pi \Rightarrow \max_n KL(\pi(\cdot; u), \pi'(\cdot; u)) = 0$

$L_n(\pi) = \eta(\pi)$ when $\pi' = \pi$

$$\underbrace{\eta(\pi)} \geq \underbrace{\max \text{ R.H.S.}} \geq \underbrace{\eta(\pi)}$$

\Rightarrow Conservative policy gradient (Kekade & Langford 2002)

\Rightarrow Trust Region Policy optimization (Schulman 2017)
(TRPO)

e.g. $\max_{\pi'} L_n(\pi')$

\rightarrow s.t. $\max_n KL(\pi(\cdot; u), \pi'(\cdot; u)) < \delta$

You can extend it to POMDP.

Linear Quadratic Regulators (LQR)

$$\rightarrow (X = \mathbb{R}^n, A = \mathbb{R}^n, A, B, R, Q, P, \Sigma)$$

$$A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times m}, \Sigma, Q \in \mathbb{R}^{n \times n}, R \in \mathbb{R}^{m \times m}$$

where R is PD, Q, Σ are PSD

$$X_1 \sim P_1$$

action, control input
↓

$$X_{t+1} = A X_t + B U_t + W_t$$

$$\Rightarrow C_t = \|X_t\|_Q^2 + \|U_t\|_R^2$$

mean zero noise with covariance Σ

Policy to choose U_t to minimize undiscarded infinite horizon cost.

\Rightarrow The optimal policy π ; $U = \underline{K} X$; $K \in \mathbb{R}^{m \times n}$ exists, and is deterministic.

$$K = - \left(R + B^T P B \right)^{-1} B^T P A$$

P is the unique PD solution to the discrete time algebraic Riccati equation (DARE)
(similar to Poisson or dynamic programming)

Monday, November 30, 2020

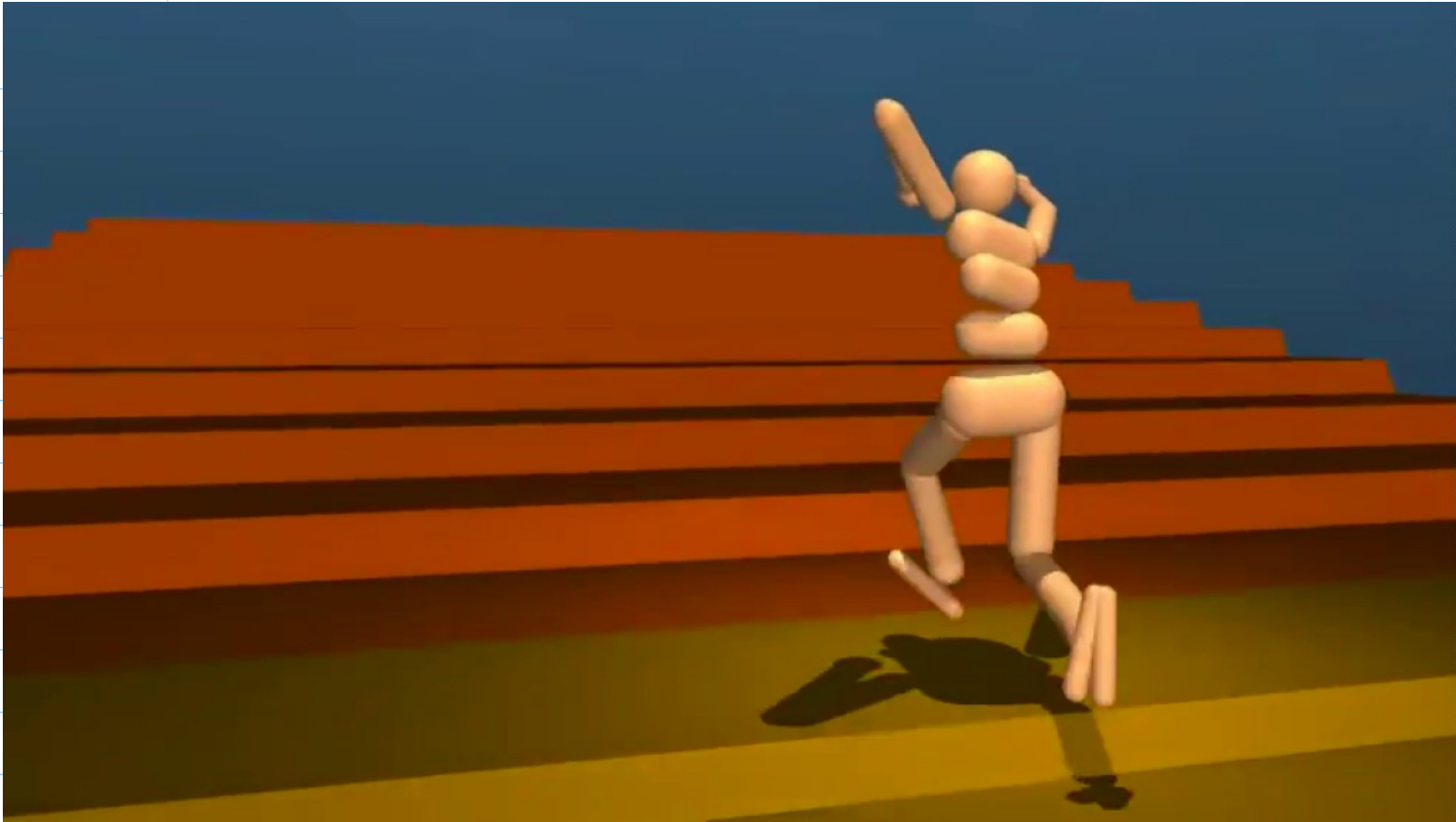
$$\rightarrow P = A^T P A + Q - A^T P B (R + B^T P B)^{-1} B^T P A$$

\Rightarrow Does policy gradient converge to the optimal K ?

yes. Under some regularity condition, Σ being PD

\Rightarrow Powerful.

Monday, November 30, 2020



Monday, November 30, 2020

Ti

DEEPMIND AI LEARNED HOW TO WALK