

policy gradients

goal: derive a model-free algorithm
that searches directly over
policy parameters

Recht 2018 Part 9

Williams 1992

Fazel, Ge, Kakade, Mesbahi 2018

Kakade 2001

Schulman, Levine, Moritz, Jordan, Abbeel 2015

Maria, Guy, Recht 2018

• suppose we apply parametric, randomized
policy $\mu: X \times \Psi \rightarrow \Delta(\mathcal{U})$ to

SDE $x^+ \sim P(x, u)$

– this determines a random process

$$u \sim \mu(x; \psi), \quad x^+ \sim P(x, u)$$

with samples / roll-outs / trajectories

$$\tau: [0, t) \rightarrow X \times \mathcal{U}$$

$$: s \mapsto (x_s, u_s)$$

- express probability distribution over trajectories i.t.o. P, μ

$$- P_{\psi}(\tau) = \prod_{s=0}^{t-1} P(x_{s+1} | x_s, u_s) \cdot \mu_{\psi}(u_s | x_s)$$

- suppose we want to choose $\psi \in \Psi$ to minimize expected finite-horizon cost,

$$E_{P_{\psi}}[c(\tau)] = E_{P_{\psi}} \left[\sum_{s=0}^t \mathcal{L}(s, x_s, u_s) \right],$$

i.e. we want to solve $\min_{\psi \in \Psi} E_{P_{\psi}}[c(\tau)]$

- let's see what happens when we apply gradient descent:

$$D_{\psi} E[c(\tau)] = D_{\psi} \int c(\tau) \cdot P_{\psi}(\tau) d\tau$$

$$= \int c(\tau) \cdot D_{\psi} P_{\psi}(\tau) d\tau \quad \text{- assuming Leibniz' rule applies}$$

$$= \int c(\tau) \cdot \left(\frac{D_{\psi} P_{\psi}(\tau)}{P_{\psi}(\tau)} \right) P_{\psi}(\tau) d\tau \quad \text{- assuming } P_{\psi}(\tau) \neq 0$$

$$= \int (c(\tau) \cdot D_{\psi} \log P_{\psi}(\tau)) \cdot P_{\psi}(\tau) d\tau \quad \text{- substitution}$$

$$= E_{P_\psi} [c(\tau) \cdot D_\psi \log P_\psi(\tau)] - \text{def. of expectation}$$

- here's the magic:

→ compute $D_\psi \log P_\psi(\tau)$, simplify

$$\begin{aligned} - \log P_\psi(\tau) &= \log \prod_{s=0}^{t-1} P(x_{s+1} | x_s, u_s) \cdot \mu_\psi(u_s | x_s) \\ &= \sum_{s=0}^t P(x_{s+1} | x_s, u_s) + \sum_{s=0}^t \mu_\psi(u_s | x_s) \end{aligned}$$

$$- \text{so } D_\psi \log P_\psi(\tau) = \sum_{s=0}^{t-1} D_\psi \mu_\psi(u_s | x_s)$$

* we only need to know derivative of the policy,
not the system model !

* furthermore, Monte Carlo gives us:

$$\begin{aligned} E_{P_\psi} [c(\tau) \cdot D_\psi \log P_\psi(\tau)] \\ \simeq \frac{1}{N} \sum_{n=1}^N c(\tau_n) \cdot D_\psi \log P_\psi(\tau_n) \end{aligned}$$

◦ if this seems like magic: it is...

- same derivation applies to any

- same derivation applies to any (even deterministic!) optimization problem:

instead of $\min_{u \in \mathcal{U}} c(u)$

consider $\min_{\mu \in \Delta(\mathcal{U})} \mathbb{E}_{\mu}[c(u)]$

where we've randomized the policy

- if we then restrict to policies parameterized by $\psi \in \Psi$ then exactly the same derivation yields

$$\nabla_{\psi} \mathbb{E}_{\psi}[c(u)] = \mathbb{E}_{\psi}[c(u) \cdot \nabla_{\psi} \log \mu_{\psi}(u)]$$

- applying Monte Carlo to approximate this gradient is termed the "REINFORCE" algorithm

NOTE: randomization and parameterization both imply any solution we obtain is sub-optimal, i.e. provides an upper bound on

... provides an upper bound on
actual minimal cost achievable

- since, at the end of the day, this is
a "gradient"-based method,

("score quotes" added because the
desired gradient does exist, but
is only ever (poorly) approximated)

there are obvious variations based
on techniques that (can) improve
on steepest descent:

- Newton - Raphson method
uses the Hessian to improve
scaling / avoid chattering:

$$\psi^+ = \psi - \left[D_\psi^2 E_\psi [c(\tau)] \right]^{-1} \cdot D_\psi E_\psi [c(\tau)]$$

termed "natural policy gradient"

Kakade 2001

- trust-region methods choose step size by solving auxiliary problem:

$$\psi^+ = \psi + \min_{\|\delta\| \leq 1} D_\psi E_\psi[c(\tau)] \cdot \delta$$

where the norm $\|\cdot\|$ chosen could be:

$\|\cdot\|_2$ (standard 2-norm): steepest descent

$\|\cdot\|_{D_\psi^2 E_\psi[c(\tau)]}$ (Hessian norm): natural P.G.

... or another of your choosing
termed "trust-region policy
optimization" (TRPO)

Schulman et al 2015

code: github.com/rll/rllab