

4 Policy Gradient Algs.

- Policy search \rightarrow instead of maintaining estimates of value funcⁿ, search in the space of policies.
 - \hookrightarrow simple, better convergence, robust to partial observability! (more popular)
 - \rightarrow Direct policy search \rightarrow e.g. genetic algos.
 - \hookrightarrow we will use policy gradient approaches.
 - \hookrightarrow most of the times description of policy π is much simpler than the value funcⁿ. Better to go for better approximator than going for wrong final values.
- What is partial observability? \rightarrow We don't have access to true state. But we get something like contextual state. Then, we add more such contextual states, to estimate the further values.
- Policy Gradient Methods
 - \hookrightarrow depends on $\theta \rightarrow$ (action preferences, mean, variance, weights of a neural network).
 - \rightarrow The idea is to modify the current policy params directly instead of

estimating the α^n values.

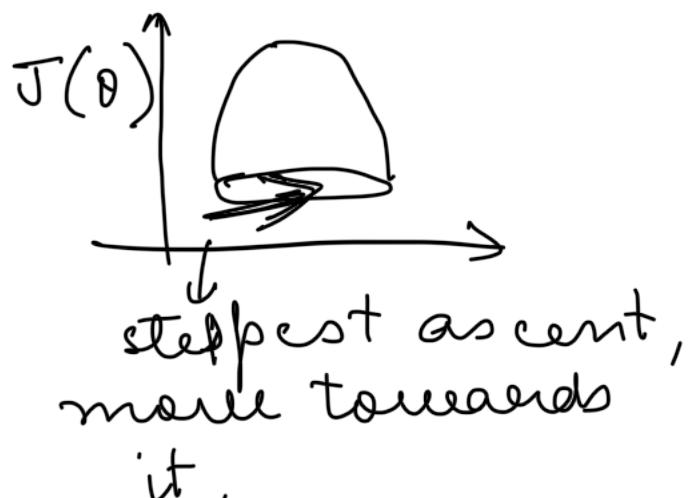
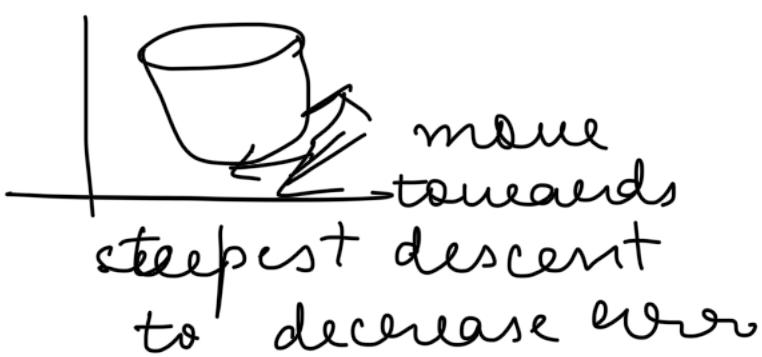
Max. $J(\theta) = E(a_t)$ → performance obj. funcⁿ for policy
 θ updaⁿ rule $\Rightarrow \theta \leftarrow \theta + \alpha \nabla J(\theta)$ $\frac{J}{\nabla}$.

$$J(\theta) = E(a_t) = E(r_t)$$

↳ expected total reward.

take as small step in the direction of the gradient

GD



- let's talk about simplified setting
↳ immediate reward or multi-arm bandit.

set of n -arms, pick the arm with highest expected reward.

$$J(\theta) = E(a_t) = \sum_a q_{\pi^*}(a) \pi_\theta(a)$$

true

expected

reward

mean of the gaussian.

- Stochastic Gradient Ascent

$$\nabla J(\theta) = \sum_a q_{\pi^*}(a) \nabla J_{\pi_\theta}(a) = \mathbb{E}_{\pi_\theta} \{q_{\pi^*}(a)\}$$

$$= \sum_a q_{\pi^*}(a) \frac{\nabla J_{\pi_\theta}(a)}{\pi_\theta(a)} \pi_\theta(a)$$

I don't know q_{π^*} . \Rightarrow sample it.

Now estimate the gradient from N samples.

$$\hat{\nabla} J(\theta) = \frac{1}{N} \sum_{i=1}^N a_i \underbrace{\frac{\nabla J_{\pi_\theta}(a_i)}{\pi_\theta(a_i)}}_{\text{cl'd the likelihood ratio.}}$$

kind of importance sampling. \leftarrow

- Reinforcement \rightarrow step sizes are small enough, eventually you will converge to true local maxima.

$$\Delta \theta_t = \alpha a_t \frac{\nabla J_{\pi_\theta}(a_t)}{\pi_\theta(a_t)}$$

$$\Delta \theta_t = \alpha a_t \frac{\partial \ln \pi_\theta(a_t)}{\partial \theta}$$

With baseline \rightarrow it helps to stabilize the learning

$$\Delta \theta_t = \alpha (a_t - b_t) \frac{\partial \ln \pi_\theta(\hat{a}_t)}{\partial \theta}$$

reinforcement baseline \rightarrow - eligibility

b_t should n't depend on a_t
 $a_t \leq "$ " " $a_t \Rightarrow$ otherwise
it leads to bias.

L2 Reinforcement

↳ Williams → 1992 with arbitrary rewards.

Eq. binary bandit problem → generalised L_{R-I}

$$\pi(\theta, a) = \begin{cases} \theta & \text{if } a=1 \\ 1-\theta & \text{if } a=0 \end{cases} \quad \frac{\partial \ln \pi}{\partial \theta} =$$

baseline $b=0$ & $\alpha = \beta \cdot \theta (1-\theta)$ $\frac{\alpha - \theta}{\theta(1-\theta)}$

$$\boxed{\Delta \theta = \beta \cdot a \cdot (\alpha - \theta)}$$

$$\begin{array}{lll} a=1 & \alpha > 0 & \theta \uparrow \\ a=0 & \alpha > 0 & \theta \downarrow \end{array}, \quad \begin{array}{lll} \alpha < 0 & \theta \downarrow \\ \alpha > 0 & \theta \uparrow \end{array}$$

$\alpha = 0 \Rightarrow$ no change

if penalty, no Δ , if +ve reward,
we react to it.

Eq. Softmax

Set a baseline to get avg. of obs. rewards.

$$b_t = \bar{u}_t = \bar{u}_{t-1} + \beta \cdot (u_t - \bar{u}_{t-1})$$

we modify the softmax action select.

$$\Delta \theta_i = \alpha \cdot (\bar{a} - a) \underbrace{(1 - \pi(\theta, a_i))}_{\rightarrow b^t}$$

$$\pi(\theta, a_i) = \frac{e^{\theta_i}}{\sum_{j=1}^n e^{\theta_j}}$$

$$\frac{\partial \ln \pi(\theta, a_i)}{\partial \theta_i} = 1 - \pi(\theta, a_i)$$

this is how we compute the characteristic eligibility for softmax action selection.

Eq. For continuous actions

↳ gaussian distribution to select actions

$$\pi(a, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(a-\mu)^2}{2\sigma^2}}$$

$$\Delta \mu = \alpha \cdot (\bar{a} - \bar{\mu}) (a - \mu)$$

$$\Delta \sigma = (\alpha/\sigma) \cdot (\bar{a} - \bar{\mu}) ((a - \mu)^2 - \sigma^2)$$

$$\alpha_\mu = \alpha_\sigma = \alpha \sigma^2$$

$$b_\mu = b_\sigma = \bar{a}$$

↳ Policy Gradient Theorem

- for episodic task, we define the performance

by assuming that every episode starts from state s_0 (non-random)

$$J(\theta) = v_{\bar{\pi}_\theta}(s_0)$$

\hookrightarrow true value funcⁿ given a policy $\bar{\pi}_\theta$

- Th. (Marbach.)

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\bar{\pi}}(s, a) \nabla \bar{\pi}(a|s, \theta)$$

It has an analytical expression for the performance gradient with respect to policy param θ .

$$\begin{aligned} \nabla v_{\bar{\pi}}(s) &= \nabla \left[\sum_a \bar{\pi}(a|s) q_{\bar{\pi}}(s, a) \right] \\ &= \sum_a \left[\nabla \bar{\pi}(a|s) q_{\bar{\pi}}(s, a) + \bar{\pi}(a|s) \nabla q_{\bar{\pi}}(s, a) \right]. \end{aligned}$$

:

$$\nabla v_{\bar{\pi}}(s) = \sum_{x \in S} \sum_{k=0}^{\infty} p_{\bar{\pi}}(s \rightarrow x, k, \bar{\pi})$$

$$\sum_a \nabla \bar{\pi}(a|x) q_{\bar{\pi}}(x, a)$$

$p_{\bar{\pi}}$ is transitioning probab. from state s to x in k steps under policy $\bar{\pi}$.

$$\nabla J(\theta) = \nabla v_{\bar{\pi}}(s_0)$$

$$\propto \sum_s \mu(s) \sum_a \nabla \bar{\pi}(a|s) q_{\bar{\pi}}(s, a)$$