

1) Q Function Approximation

$$\{(s_i, a_i, y_i)\}_{i=1}^N$$

rollout

MC

TD

Bellman opt.

$$\hat{Q} = \arg \min_{\theta} \sum_{i=1}^N \underbrace{(\underbrace{Q_{\theta}(s_i, a_i)} - y_i)^2}_{\text{error}}$$

$$\mathcal{Q} = \{Q_{\theta} : \theta \in \mathbb{R}^d\}$$

Incremental updates: gradient-based method

$$\nabla_{\theta} [(Q_{\theta}(s_i, a_i) - y_i)^2] = 2 \underbrace{(Q_{\theta}(s_i, a_i) - y_i)}_{\text{error}} \underbrace{\nabla_{\theta} Q_{\theta}(s_i, a_i)}$$

update

$$\theta \leftarrow \begin{cases} \theta - \alpha (Q_{\theta}(s_i, a_i) - y_i) \nabla_{\theta} Q_{\theta}(s_i, a_i) \\ y_i = r_i + \gamma \underbrace{Q_{\theta}(s_{i+1}, a_{i+1})} \end{cases}$$

How to choose (s_i, a_i, y_i) to update with respect to:

1) online GD uses (s_t, a_t, y_t)

2) "experience replay": store incoming data
resample $\{(s_i, a_i, y_i)\}_{i=1}^N$ randomly from stored data
for $i=1, \dots, N$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(Q_{\theta}(s_i, a_i) - y_i)^2$$

2) Optimization & Gradient Ascent

Goal: find π^* (approximately). Why bother with α^* ?

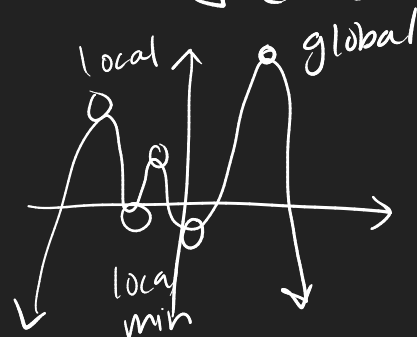
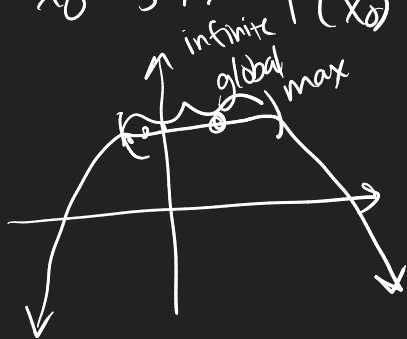
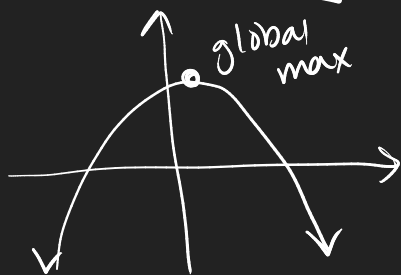
suppose parametrized policy π_θ

objective function: $J(\theta) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid P, \pi_\theta, s_0 \right]$

Maxima & Minima: $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}$

A global max is x_0 s.t. $f(x_0) \geq f(x) \quad \forall x \in \mathbb{R}^d$

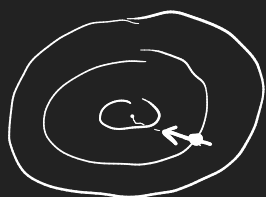
A local max is x_0 s.t. $f(x_0) \geq f(x) \quad \forall \|x - x_0\| \leq \epsilon$
 $\exists \epsilon > 0$



An ascent direction is $v \in \mathbb{R}^d$ s.t.
 $f(x + \alpha v) > f(x)$ for some $\alpha > 0$



If f is differentiable, the $\nabla f(x)$ is the steepest ascent direction



Gradient Ascent:

initial x_0

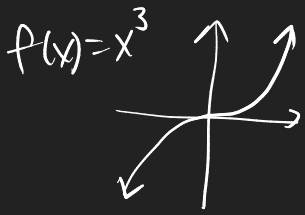
for $t=0, 1, \dots$

$$x_{t+1} = x_t + \alpha \nabla f(x_t)$$

↑ step size

$$\underbrace{f(x)}_{\max f(x)} \approx \underbrace{f(x_t) + \nabla f(x_t)^T (x - x_t)}_{x - x_t \propto \nabla f(x_t)}, \text{ when } x \& x_t \text{ close} \\ \propto \text{keeping close}$$

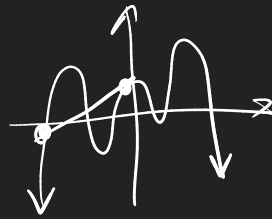
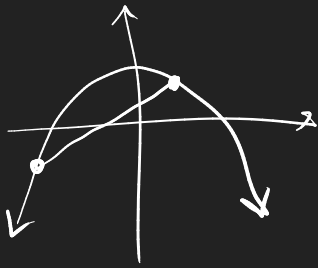
Fact $\nabla f(x_0) = 0 \Leftrightarrow x_0$ local max



Critical Points: $\forall x$ s.t. $\nabla f(x) = 0$

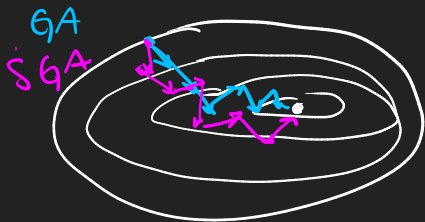
- local/global max/min
- saddle points

Concave function has critical point(s) which are global maxima



3) Stochastic Gradient Ascent

Instead $\nabla f(x_t)$, we have g_t s.t. $\mathbb{E}[g_t] = \nabla f(x_t)$



Alg: SGA

init x_0
for $t=0, 1, \dots$

$$x_{t+1} = x_t + \alpha g_t$$

ex: SGD for ERM

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i)$$

$$g_t = \nabla_{\theta_t} \ell(f_{\theta_t}(x_i), y_i)$$

$x_i, y_i \sim$ uniformly sampled

$$\mathbb{E}[g_t] = \sum_{i=1}^n \nabla_{\theta_t} \ell(f_{\theta_t}(x_i), y_i) \cdot \frac{1}{N}$$

Theorem: Assume:

-1) $f(x)$ is β smooth $\|\nabla f(x) - \nabla f(x')\|_2 \leq \beta \|x - x'\|_2$

-2) $\sup_x |f(x)| \leq M$

[3) $\mathbb{E}[g(x)] = \nabla f(x)$

[4) $\mathbb{E}[\|g(x)\|_2^2] \leq \sigma^2$

Then SGA with $g_t = g(x_t)$ converges

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla f(x_t)\|_2\right] \leq \sqrt{\frac{\beta \sigma^2 M}{T}}$$

$$\alpha = \sqrt{\frac{M}{\beta \sigma^2 T}}$$

Question can a sampled trajectory $\tau = (s_0, a_0, \dots)$ directly give us estimate of $\nabla J(\theta)$?

A: no! not knowing transition P or reward r is like not knowing loss function in SL

$$J(\theta) = \mathbb{E}_{w \sim N(0,1)} [s_1^2 \mid s_1 = f(s_0, a_0, w), a_0 = \pi_\theta(s_0)]$$

$$\nabla J(\theta) = \nabla_\theta \mathbb{E}_w [f(s_0, \pi_\theta(s_0), w)^2]$$

$$= \mathbb{E}_w [\nabla_\theta f(s_0, \pi_\theta(s_0), w)^2] \neq \mathbb{E}[s_1^2]$$

one trajectory: (s_0, a_0, s_1)

Next Lecture: Derivative-Free Optimization

can't access $\nabla f(x)$ or even estimate

but can access $f(x)$