· Recall:

Discounted state visitation:

$$d_{s_0}^{\pi}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr[S_t = s | S_0]$$

• To sample from this distribution.

  Start at $s_0$ & simulate $\pi$; Accept a state with prob $1-\gamma$,

  prob of accepting $s$ at time $t$?
  don't accept at $t=0, t=1, \ldots, t=t-1$,
  accept at $t = (1-\gamma)\gamma^t \Pr[S_t = s | S_0]$.

  ∴ the state $s$ is distributed as
  $$d_{s_0}^{\pi}(s) !$$

If $\tau$ is a trajectory,

the unconditional distribution $\mathbb{P}^{\pi}_{\mu}(\tau)$

under $\pi$ starting with initial distribution $\mu$ is

$$\mu(s_0)\,\pi(a_0|s_0)\,P(s_1|s_0\,a_0)\cdots$$

Notation: $\quad d^{\pi}_{\mu}(s) = \underset{s_0 \sim \mu}{\mathbb{E}}\left[d^{\pi}_{s_0}(s)\right].$

Given $f: S \times A \to \mathbb{R}$;

$$\underset{\tau \sim \mathbb{P}^{\pi}_{\sigma}}{\mathbb{E}}\left[\sum_{t=0}^{\infty}\gamma^t f(s_t, a_t)\right]$$

<span style="color:red">Expected discounted value of $f$ along the trajectory</span>

$$\overset{?}{=} \frac{1}{1-\gamma}\underset{s \sim d^{\pi_\theta}}{\mathbb{E}}\;\underset{a \sim \pi_\theta(\cdot|s)}{\mathbb{E}}\; f(a,s).$$

In terms of sampling,

$$\mathbb{E}_{\tau \sim \mathbb{P}_r^{\pi}} \left[ \sum \gamma^t f(s_t, a_t) \right] \quad \text{is:}$$

Sample a trajectory; compute the
discounted value of $f$ over the trajectory;
Find the expected value;

- Run the markov chein: with prob $1-\gamma$
  select state, and then on action $\left( \pi(\cdot | s) \right)$
  and compute $\dfrac{f(s,a)}{1-\gamma}$

prob of picking $s$ in time $t$ ?
$$= (1-\gamma) \gamma^t \Pr\left[ s_t = s | s_0 \right]$$

Pro of picking action $a$ $\pi(\cdot | s)$

$$\therefore \text{get } \dfrac{(1-\gamma) \gamma^t \boxed{\Pr\left[ s_t = s | s_0 \right] \cdot \pi(\cdot | s)} f(s,a)}{1-\gamma}$$

Regroup this sum and get

$$\mathbb{E}_{\tau \sim \mathbb{P}_\delta^\pi} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right]$$

$$\text{||}$$

$$\therefore \quad \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_\delta}} \mathbb{E}_{a \sim \pi(\cdot|s)} f(s,a)$$

$$=$$

<u>Lemma</u>: For all policies $\pi, \pi'$, and initial

dist $\mu$,

$$V^\pi(\mu) - V^{\pi'}(\mu) = \mathbb{E}_{\tau \sim \mathbb{P}_\mu^\pi} \left[ \sum \gamma^t A^{\pi'}(s_t, a_t) \right]$$

advantage of a policy;  $A^\pi(s,a) \quad Q^\pi(s,a) - V^\pi(s)$

For a fixed state $s$,

$\mathbb{P}_s^{\pi}$ denote distribution over trajectories with $s_0 = s$;

$V^{\pi}(s) - V^{\pi'}(s)$

$= (1-\gamma) \underset{\tau \sim \mathbb{P}_s^{\pi}}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] - V^{\pi'}(s)$

$= \underset{\tau \sim \mathbb{P}_s^{\pi}}{\mathbb{E}} \left[ (1-\gamma)\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) + \sum \gamma^t V^{\pi'}(s_t) - \sum \gamma^t V^{\pi'}(s_t) \right]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad - V^{\pi'}(s)$

$= \underset{\tau \sim \mathbb{P}_s^{\pi}}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t \left[ (1-\gamma) r(s_t, a_t) + \gamma V^{\pi'}(s_{t+1}) - V^{\pi'}(s_t) \right] \right]$

$= \underset{\tau \sim \mathbb{P}_s^{\pi}}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t \left[ (1-\gamma) r(s_t, a_t) + \gamma \mathbb{E}\left[ V^{\pi'}(s_{t+1}) \mid s_t, a_t \right] \right. \right.$
$\qquad\qquad\qquad\qquad \left. \left. - V^{\pi'}(s_t) \right] \right]$

$\left( \underset{s_t, a_t, s_{t+1} \cdots}{\mathbb{E}} V^{\pi}(s_{t+1}) = \underset{s_t, a_t}{\mathbb{E}} \underset{s_{t+1}}{\mathbb{E}}\left[ V^{\pi}(s_{t+1}) \mid s_t, a_t \right] \right)$

$$= \underset{\tau \sim \mathbb{R}^{\pi}}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t \left( Q^{\pi'}(s_t, a_t) - V^{\pi'}(s_t) \right) \right]$$

$$= \underset{\tau \sim \mathbb{R}^{\pi}}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t \, A^{\pi'}(s_t, a_t) \right],$$

## Policy gradient:

Discounted total reward of a trajectory.

$$R(\tau) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

Now: $\quad V^{\pi_\theta}(\mu) = \underset{\tau \sim \mathbb{R}^{\pi_\theta}_\mu}{\mathbb{E}} \left[ R(\tau) \right]$

$=$

THM: POLICY GRADIENT'S:

$$\nabla V^{\pi_\theta}(\mu) = \mathop{\mathbb{E}}_{\tau \sim \mathbb{P}_\mu^{\pi_\theta}} \left[ R(\tau) \sum_{t=0}^{\infty} \nabla \log \pi_\theta (a_t | s_t) \right]$$

↑ ⌐ In terms of discounted reward $\beta$

## REINFORCE:

$$\nabla V^{\pi_\theta}(\mu) = \nabla \mathop{\mathbb{E}}_{\tau \sim \mathbb{P}_\mu^{\pi_\theta}} (R(\tau))$$

$\underline{Q}: \int \quad ?$

Countability issues.

$$= \nabla \sum_\tau R(\tau) \, \mathbb{P}_\mu^{\pi_\theta}(\tau)$$

$$= \sum_\tau R(\tau) \, \nabla \mathbb{P}_\mu^{\pi_\theta}(\tau)$$

$$= \sum_\tau R(\tau) \, \mathbb{P}_\mu^{\pi_\theta}(\tau) \, \nabla \log \mathbb{P}_\mu^{\pi_\theta}(\tau)$$

$$= \sum_\tau R(\tau) \, \mathbb{P}_\mu^{\pi_\theta}(\tau) \, \nabla \log \left[ \mu(s_0) \pi_\theta(a_0|s_0) P(s_0, a_0, s_1) \cdots \right]$$

$$= \sum_\tau R(\tau) \, \mathbb{P}_\mu^{\pi_\theta}(\tau) \sum_{t=0}^{\infty} \nabla \log \pi_\theta (a_t | s_t)$$

$$= \mathop{\mathbb{E}}_{\tau \sim \mathbb{P}_\mu^{\pi_\theta}} R(\tau) \sum_{t=0}^{\infty} \nabla \log \pi_\theta (a_t | s_t)$$

$$\nabla V^{\pi_\theta}(\mu) = \underbrace{\frac{1}{1-\gamma} \underset{s\sim d^{\pi_\theta}}{\mathbb{E}} \underset{a\sim\pi(\cdot|s)}{\mathbb{E}} \left[ Q^{\pi_\theta}(s,a) \nabla \log \pi_\theta(a|s) \right]}_{}$$

$$\nabla V^{\pi_\theta}(\mu) =$$
$$\underset{\tau\sim P^{\pi_\theta}_\mu}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t \underset{\uparrow}{Q^{\pi_\theta}(s_t, a_t)} \nabla \log \pi_\theta(s_t|a_t) \right]$$

In terms of action value

<span style="color:red">An unbiased estimate!</span>

$$\cdot \; \nabla V^{\pi_\theta}(\mu) = \frac{1}{1-\gamma} \underset{s\sim d^{\pi_\theta}}{\mathbb{E}} \underset{a\sim\pi_\theta(\cdot|s)}{\mathbb{E}} \left[ A^{\pi_\theta}(s,a) \nabla \log \pi_\theta(a|s) \right]$$

In terms of advantage

$$= \underline{\text{Gradient ascent}:}$$

$$f: \mathbb{R}^d \to \mathbb{R} \; \text{is} \; \beta \; \text{smooth if}$$
$$\| \nabla f(\omega) - \nabla f(\omega') \| \leq \beta \| \omega - \omega' \|$$

Update: $\theta_{t+1} = \theta_t + \eta \nabla V^{\pi_\theta}(\mu)$.

Lemma: Assume $V^{\pi_\theta}$ is $\beta$

Smooth $\forall \theta$, Assume $V^{\pi_\theta}$ is bounded below

by $V_*$, Using $\eta = 1/\beta$, $\forall T$

$$\min_{t \leq T} \left\| \nabla V^{(t)}(\mu) \right\| \leq \frac{2\beta \left( V^*(\mu) - V^0(\mu) \right)}{T}.$$

when

$$\boxed{T \geq \frac{1}{\varepsilon} 2\beta \left[ V^*(\mu) - V^0(\mu) \right] ,}$$ , one

of the gradients is $\leq \varepsilon$, (ie) close to a

stationary point then!

Unbiased estimate of gradients:

Assume we have sampler, $\tau \sim \mathbb{P}_r^{\pi_\theta}$;

For a trajectory $\tau$ define: <span style="color:red">Discounted Gain from state $t$ onwards</span>

$$\widehat{Q^{\pi_\theta}}(s_t, a_t) := (1-\gamma) \sum_{t'=t} \gamma^{t'-t} r(s_t, a_{t'})$$

$$\widehat{\nabla V^{\pi_\theta}(\mu)} := \sum_{t=0}^{\infty} \gamma^t \widehat{Q^{\pi_\theta}}(s_t, a_t) \nabla \log \pi_\theta(a_t|s_t)$$

Claim: $\underset{\tau \sim \mathbb{P}_r^{\pi_\theta}}{\mathbb{E}}\left[ \widehat{\nabla V^{\pi_\theta}(\mu)} \right] = \nabla V^{\pi_\theta}(\mu)$

$\uparrow$ unbiased estimator;

Proof:

$$\underset{\tau \sim \mathbb{P}_r^{\pi_\theta}}{\mathbb{E}}\left[ \sum_{t=0}^{\infty} \gamma^t \widehat{Q^{\pi_\theta}}(s_t, a_t) \nabla \log \pi_\theta(a_t|s_t) \right]$$

$$= \underset{\tau \sim \mathbb{P}_r^{\pi_\theta}}{\mathbb{E}}\left[ \sum_{t=0}^{\infty} \gamma^t \underbrace{\mathbb{E}\left[ \widehat{Q^{\pi_\theta}}(s_t, a_t) \mid s_t, a_t \right]}_{= Q(s_t, a_t)} \nabla \log (\,) \right]$$

$$= \underset{\tau \sim \mathbb{P}_r^{\pi_\theta}}{\mathbb{E}}\left[ \sum_{t=0}^{\infty} \gamma^t Q(s_t, a_t) \nabla \log(a_t|s_t) \right] = \nabla V^{\pi_\theta}(\mu)$$

Note $\widehat{Q^{\pi_\theta}}(s_t, a_t)$ is an unbiased estimator

of $Q(s_t, a_t)$.

Gives us:

Initialize $\theta_0$:

2 For $t = 0, 1, \cdots$

   (a) Sample $\tau \sim P_\mu^{\pi_\theta}$

   (b) $\theta_{t+1} = \theta_t + \eta_t \widehat{\nabla V^{\pi_\theta}}(\mu)$

- To get $\widehat{\nabla V^{\pi_\theta}}(\mu)$ from the sample

Compute $\widehat{Q^{\pi_\theta}}(s_z, a_t)$ for $t$, and use in

$\widehat{\nabla^{\pi_\theta}}(\mu)$.

Ignore that $\tau$ is $\infty$; Truncate

Algorithm: Stochastic Gradient Ascent on $J$.
REINFORCE

- Initialize $\theta$ arbitrarily.

- for each episode do

    Generate $S_0 A_0 R_0 S_1 A_1 \cdots S_{L-1} A_{L-1} R_{L-1}$

    using $\theta$.

    For each $t$ compute $\widehat{G_t} = \overbrace{Q^{\pi_\theta}(S_t, a_t)}$

    $$\widehat{\nabla J(\theta)} = \sum_{t=0}^{L-1} \gamma^t \left(\boxed{G_t}\right) \frac{\partial \ln \pi (S_t | A_t ; \theta)}{\partial \theta}$$

    $$\theta = \theta + \alpha \overbrace{\nabla J(\theta)}$$

end

- Evidently $\gamma^t$ is ignored in practice
- Maybe this is stochastic gradient for a different obj

The estimate of $\widehat{Q^{\pi\theta}}(s, a)$ has high variance.

Estimate $E(x)$.

- a single sample $x_0 \rightarrow$ estimate: $x_0$.

Problem: high variance;

Suppose we take a sample of another $Y$, whose expectation $E(Y)$ we know.

Try $\hat{\mu}^a = X - Y + E(Y)$;

$$Var(\hat{\mu}^a) = Var(X - Y) = Var(X) + Var(Y) - Cov(X, Y)$$

If $Var(Y) < Cov(X, Y)$ — in business

- $Y$ is called a control variate

Want: $X, Y$ positively correlated;

## FOR REINFORCE:

$$\theta \leftarrow \theta + \alpha \sum_{t=0}^{L-1} \gamma^t \left( \overset{\pi_\theta}{Q}(s_t, a_t) - f(s_t) \right) \nabla \log \pi_\theta (a_t | s_t)$$

Here:

$f: S \to \mathbb{R}$ is a function indep $g$

$\tau$ a traj

For any function $g(s)$,

$$\mathbb{E} \left[ \nabla \log \pi(a|s) \, g(s) \right]$$

$$= \sum_a \pi(a|s) \, \nabla \log (\pi(a|s)) \, g(s)$$

$$= \sum_a \frac{\pi(a|s)}{\pi(a|s)} \nabla \pi(a|s) \, g(s)$$

$$= g(s) \sum_a \nabla \pi(a|s) = g(s) \, \nabla \sum_a \pi(a|s)$$

$$= g(s) \, \nabla (1) = \underline{\underline{0}}.$$

Now if $f(\cdot)$ is independent of $\tau$, $\forall t$.

$$\mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t f(s_t) \nabla \log \pi_\theta(a_t | s_t)\right] = 0$$

$$\sum_{t=0}^{\infty} \gamma^t \underbrace{\mathbb{E}\left[f(s_t) \nabla \log \pi_\theta(a_t | s_t)\right]}_{0}$$

$\overset{\shortmid\shortmid}{0}$

- $f(s)$ called baseline at state $s$.

What do we use for $f$?

Common choice is $v^\theta$;

- More MC algorithm for $\nu^{\pi}(\cdot)$

- $\Pr_{\gamma}^{\pi} \leftarrow$ the distribution of trajectories
  starting at $s_0$.

Recall: Generate episodes using $\pi$;

for every state $s$, $t_s$ be the first time $s$
appears in episode;

$$G_s \leftarrow \sum_{k=0}^{\infty} \gamma^k r\left(s_{t_s+k}, a_{t_s+k}\right) \quad \text{discounted}$$

$\hookrightarrow$ reward from then
on

Append $G_s$ to returns $(s)$;

For each $s$, return avg (returns $(s)$).

# Gradient based MC:

update at time t.

$$v(S_t) \leftarrow v(S_t) + \alpha \left( G_t - V(S_t) \right).$$

- This is minimizing mean squared value error;

$$\frac{1}{2} \underset{s}{\mathbb{E}} \left[ \left( v^{\pi}(s) - v(s) \right)^2 \right].$$

  - Move against the gradient of the above loss.

$$v \leftarrow v - \alpha \cdot \frac{\partial}{\partial v} \mathbb{E}_s \left[ (\ )^2 \right]$$

$$\leftarrow v - \alpha \cdot \mathbb{E} \left( v^{\pi}(s) - v(s) \right) (-1) \frac{\partial v(s)}{\partial v}$$

$$= v + \alpha \mathbb{E} \left[ \left( v(s) - v^{\pi}(s) \right) \frac{\partial v(s)}{\partial v} \right]$$

. Don't know $v^{\pi}(s)$.

So use the discounted reward from first visit to state $s$ in the sample trajectory.

Unbiased estimator of $v^{\pi}(s)$;

- **Temporal difference learning:**

  - policy evaluation algorithm

  - We don't know $P$ & $R$;

  - Get samples & learn from experience (as in MC algorithms)

  - In TD we estimate $V^{\pi}$;

**TD update:**

- if in state $s$, we take action $a$, go to $s'$ & get reward $r$

$$v(s) \longleftarrow v(s) + \alpha\left(r + v(s') - v(s)\right)$$