Define

$S_0 =$ initial state (could use $\mu(s) =$ initial distribution)

$\pi_\theta =$ Parametric policy (e.g., neural net). Will write $\pi$ sometimes for simplicity.

$\tau =$ A trajectory $\tau = (s_0, a_0, s_1, a_1, \ldots, s_H, a_H)$

Note: we will truncate at $H$ here. not sure if that matters in general

$R(\tau) =$ Return along trajectory $\tau = \sum_{t=0}^{H} r_t s_t$, wh $s_t$ is reward at time $t$.

$J(\theta) := \mathbb{E}_{\pi_\theta}(R(\tau)) = \sum_\tau \underbrace{P(\tau; \theta)}R(\tau)$

Probability of trajectory $\tau$

Observe

$$\nabla_\theta J(\theta) = \sum_\tau \nabla_\theta P(\tau; \theta) R(\tau)$$

Want to write as expectation so we can use WLLN / Sampling

$$\longrightarrow = \sum_\tau P(\tau; \theta) \frac{\nabla_\theta P(\tau; \theta)}{P(\tau; \theta)} R(\tau)$$

$$= \mathbb{E}_{\pi_\theta}\left( \nabla_\theta \ln[P(\tau; \theta)] R(\tau) \right)$$

Hence,

$$\hat{g} := \frac{1}{m} \sum_{i=1}^{m} \nabla_\theta \ln[P(\tau^{(i)}; \theta)] R(\tau^{(i)})$$

is an unbiased estimate of $\nabla_\theta J(\theta)$

Lets look more carefully at $\nabla_\theta \ln P(\tau;\theta)$

$$\nabla_\theta \ln P(\tau;\theta) = \nabla_\theta \ln \left[ \prod_{t=0}^{H} P(s_{t+1} \mid s_t, a_t) \pi_\theta(a_t \mid s_t) \right]$$

$$= \nabla_\theta \sum_{t=0}^{H} \ln P(s_{t+1} \mid s_t, a_t) + \nabla_\theta \sum_{t=0}^{H} \ln \pi_\theta(a_t \mid s_t)$$

We can compute this for a specified parametric policy class $\longrightarrow$
$$= \sum_{t=0}^{H} \nabla_\theta \ln \pi_\theta(a_t \mid s_t)$$

$$= \sum_{t=0}^{H} \underbrace{\frac{1}{\pi_\theta(a_t \mid s_t)}}_{\substack{\text{"Correction" or} \\ \text{Normalization for} \\ \text{low probability event?} \\ \text{This part is confusing}}} \underbrace{\nabla_\theta \pi_\theta(a_t \mid s_t)}_{\substack{\text{This is a vector points in direction that} \\ \text{causes greatest increase for probability of} \\ \text{playing action } a_t \text{ from state } s_t}}$$

Recall:

$$\hat{g} = \frac{1}{m} \sum_{i=1}^{m} \left( \sum_{t=0}^{H} \frac{1}{\pi_\theta(a_t \mid s_t)} \nabla_\theta \pi(a_t \mid s_t) \right) R(\tau)$$

If a Reward for a trajectory is large, then you increase the probability of repeating that trajectory.

In practice, no one does exactly this. Variance is too high.

We may exploit temporal structure of MDP to reduce the variance of the estimator. Let

$$\tau_t = \{s_0, a_0, \dots, s_t, a_t\}$$

be a truncate of $\tau$. Note

$$J(\theta) = \mathbb{E}_\tau(R) = \mathbb{E}_\tau \left( \sum_{t=0}^{H} r_t \right)$$

$$= \sum_\tau \left( \sum_{t=0}^{H} r_t \right) P(\tau; \theta)$$

$$= \sum_\tau \sum_{t=0}^{H} r_t \, P(\tau; \theta)$$

Lazy expansion, but should exploit fact that $r_t$ is independent of $\{s_{t+1}, a_{t+1}, \dots\}$

$$= \sum_\tau \sum_{t=0}^{H} r_t \, P(\tau_t; \theta)$$

$$= \sum_{t=0}^{H} \mathbb{E}_\tau(r_t)$$

$$= \sum_{t=0}^{H} \sum_\tau r_t(\tau) P(\tau_t; \theta)$$

$$= \sum_\tau \sum_{t=0}^{H} r_t(\tau) P(\tau_t; \theta)$$

Hence,

$$\nabla_\theta J(\theta) = \sum_\tau \sum_{t=0}^{H} r_t(\tau) \nabla P(\tau_t; \theta)$$

$$= \sum_\tau \sum_{t=0}^{H} r_t(\tau) \, P(\tau_t; \theta) \nabla \ln P(\tau_t; \theta)$$

$$= \sum_\tau \sum_{t=0}^{H} P(\tau_t; \theta) \, r_t(\tau) \sum_{k=0}^{t} \nabla \ln \pi_\theta(a_k | s_k)$$

considering the $t$ by pun outside sums

$$= \mathbb{E} \left( \sum_{t=0}^{H} r_t(\tau) \sum_{k=0}^{t} \nabla \ln \pi_\theta(a_k | s_k) \right)$$

$$= \mathbb{E} \left( \sum_{t=0}^{H} \sum_{k=0}^{t} r_t \nabla \ln \pi_\theta(a_k | s_k) \right) \qquad (\ast)$$
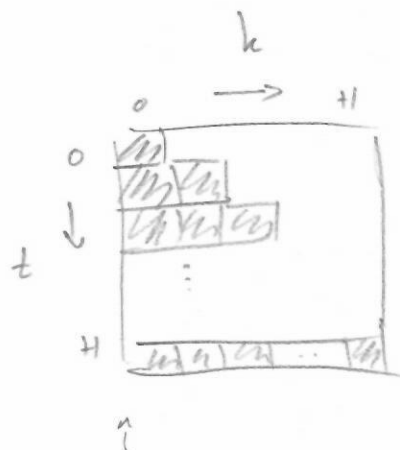
Note:

Suppose you want to sum

$$(**) \quad \sum_{t=0}^{H} \sum_{h=0}^{t} a_t b_h$$

Instead of summing rows,

Sum columns:

$$(**) = \sum_{h=0}^{H} \sum_{t=h}^{H} a_t b_h$$



Hence, $(*)$

$$= \mathbb{E} \left( \sum_{h=0}^{H} \sum_{t=h}^{H} r_t \nabla \ln \pi_\theta (a_h | s_h) \right)$$

$$= \mathbb{E} \left( \sum_{h=0}^{H} \nabla_\theta \ln \pi_\theta (a_h | s_h) \underbrace{\sum_{t=h}^{H} r_t}_{=: G_h \,=\, \text{return after time } h} \right)$$

$$= \mathbb{E} \left( \sum_{h=0}^{H} \nabla_\theta \ln \pi_\theta (a_h | s_h) \, G_h \right)$$

REINFORCE algorithm (Williams, '92)

Initialize: $\theta$

Loop forever

   Generate episode following $\pi_\theta$: $s_0, a_0, r_0, s_1, a_1, r_1, \ldots s_H, a_H, r_H$

     Loop for each step of the episode $t = 0, 1, \ldots, H$

$$G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} r_k$$

$$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla \ln \pi (a_t | s_t, \theta)$$

<u>Note:</u> It doesn't totally make sense why we're not summing up all the terms $t = 0, \ldots, H$ first, & then update $\theta$. That's the estimate, right? we really derived a different algorithm, that I'm not sure what pros & cons are. But, perspective is valuable. And will derive REINFORCE for real, next.

An alternate perspective   The Policy Gradient thrm

Defn

$q_\pi(s,a)$ = Value of taking action $a$ from State $s$ & following $\pi$ thereafter

Thrm   Let

$$J(\Theta) = v_\pi(s_0) \quad \text{Then}$$

$$\nabla J(\Theta) \propto \sum_s \mu(s) \sum_a q_\pi(s,a) \nabla \pi(a|s,\Theta)$$

pf   (Next page)

$$\nabla v_\pi(s) = \nabla \sum_a \pi(a|s) q_\pi(s,a) \qquad \forall s$$

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

$$= \sum_a \nabla \pi(a|s) q_\pi(s,a) + \pi(a,s) \nabla \sum_{s'} P(s'|s,a)(r + v_\pi(s'))$$

$$= \sum_a \nabla \pi(a|s) q_\pi(s,a) + \pi(a,s) \sum_{s'} P(s'|s,a)(\underbrace{r}_{\nabla r = 0} + \underbrace{\nabla v_\pi(s')}_{\substack{\text{Unroll this in} \\ \text{same way}}})$$

$$= \underbrace{\sum_a \nabla \pi(a|s) q_\pi(s,a)}_{①} + \pi(a,s) \sum_{s'} P(s'|s,a)$$
$$\left[ \underbrace{\sum_a \nabla \pi(a|s') q_\pi(s',a)}_{②} + \pi(a,s') \sum_{s''} P(s''|s',a) \nabla v_\pi(s'') \right]$$

Note:
- we start on std $s$
- ② is multiplied by prob of traveling to $s'$ in each $k$ time steps
- If we keep unrolling, the pattern continues

$$= \sum_{x \in S} \underbrace{\sum_{h=0}^{\infty} P(s \to x, h, \pi)}_{\text{Prob of traveling from } s \text{ to } x \text{ in exactly } k \text{ time step under } \pi} \sum_a \nabla \pi(a|x) q_\pi(x,a)$$

Hence,

$$\nabla J(\theta) = \nabla v_\pi(s_0)$$

$$= \sum_{x \in S} \left( \underbrace{\sum_{h \geq 0} P(s \to x, h, \pi)}_{\eta(s) \approx \eta(s)} \right) \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

$$= \sum_s \left( \sum_{s'} \eta(s') \right) \frac{\eta(s)}{\sum_{s'} \eta(s') \sum_a \nabla \pi(a|s) q_\pi(s,a)}, \quad \text{whr } \eta(s) = \# \text{ time steps spent on any in std } s \text{ in single episode}$$

$$= \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

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

How to use this result?

Note, thrm implies

$$\nabla J(\theta) \propto \mathbb{E}_{s \sim \mu}\left[ \sum_a q_\pi(s,a) \nabla \pi(a|s,\theta) \right]$$

We know that $s_t \sim \mu$ as $t$ large. So, for $t$ large this is same as

$$\approx \mathbb{E}_\pi\left[ \sum_a q_\pi(s,a) \nabla \pi(a|s,\theta) \right]$$

Could use estimator

$$\hat{g} = \sum_a \hat{q}(s_t,a,w) \nabla \pi(a|s_t,\theta) \qquad \text{for } s_t \text{ generated by } \pi.$$

But, this requires $\hat{q}$. Can we replace sum over actions w/ an expectation Somehow & do MC trick? Backing up, again we have

$$\nabla J(\theta) \propto \mathbb{E}_{s \sim \mu}\left[ \underbrace{\sum_a q_\pi(s,a) \nabla \pi(a|s,\theta)}_{} \right]$$

Want to write this as $\mathbb{E}$

$$= \mathbb{E}_{s \sim \mu}\left[ \sum_a \pi(a|s,\theta) q_\pi(s,a) \frac{\nabla \pi(a|s,\theta)}{\pi(a|s,\theta)} \right]$$

$$= \mathbb{E}_{s \sim \mu}\left[ \mathbb{E}_{a \sim \pi}\left[ q(s,a) \frac{\nabla \pi(a|s,\theta)}{\pi(a|s,\theta)} \right] \right]$$

Estimator: Let $s_0, a_0, \dots$ be trajectory generated by $\pi$ & let $t = $ "large"

Let

$$\hat{g} = q(s_t,a_t) \frac{\nabla \pi(a_t|s_t,\theta)}{\pi(a_r|s_t,\theta)} = \overset{\overset{\text{should be}}{\text{return}}}{G_t} \frac{\nabla \pi(a_t|s_t,\theta)}{\pi(a_t|s_t,\theta)}$$

This yields the exact same REINFORCE algorithm from before.

- Note: In ~~[faded]~~
  - Differences from prev. derivation
  - Big assumption that everyone completely ignores/glosses over:
    Need $t$ to be large for $S_t \overset{approx}{\sim} \mu$. Yet, we sample only the beginning of a trajectory, not its tail.

## REINFORCE w/ Baseline

Let $b(s)$ be some arbitrary function depending only on state $s$.

Claim:

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a \left( q_\pi(s,a) - b(s) \right) \nabla \pi(a|s, \theta)$$

This follows from the policy gradient theorem and the fact that we're actually subtracting zero

$$\sum_a b(s) \nabla \pi(a|s, \theta) = b(s) \nabla \sum_a \pi(a|s, \theta) = b(s) \nabla 1 = 0$$

A good choice of baseline can yield a lower variance unbiased estimator of $\nabla J(\theta)$. Standard choice: $b(s) \approx V(s) =$ Value S

# REINFORCE w/ baseline :

Input: parametrize $\pi(a|s,\theta)$, $\hat{v}(s;w)$
Algorith Params: step sizes $\alpha^w$, $\alpha^\theta$
Initalize $\theta$ & $w$.

Loop forever:

    Generate episode: $s_0, a_0, r_1 s_1, a_1, r_2, \dots s_H, a_H, r_H$

    Loop for each step of episode $t=0,\dots, H$

$$G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} r_k$$

$$\delta \leftarrow G - \hat{v}(s_t, w)$$

$(\star)$    $w \leftarrow w + \alpha^w \delta \nabla \hat{v}(s_t, w)$

$$\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla \ln \pi(a_t | s_t, \theta)$$

Question: why do we modulate $\hat{v}$ update by $\delta$ in $(\star)$?

    Ans: We're doing least squares fitting. Want to take a step
in $\nabla$ of $\frac{1}{2}(\hat{v}(s,w) - G)^2$

$$\nabla_w \tfrac{1}{2}(\hat{v}(s,w) - G)^2 = (\hat{v} - G)\nabla_w \hat{v}(s,w) = \delta \nabla_w \hat{v}(s,w)$$