RECALL: REINFORCE WITH BASELINE.

Want: $X, Y$ positively correlated;

FOR REINFORCE:

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

Base line

Here:

$$f: S \rightarrow \mathbb{R} \quad \text{is a function indep } g$$

$\tau$ a traj

For any function $g(s)$,

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

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

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

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

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

Recall - Reinforce algorithm.

Algorithm: Stochastic Gradient Ascent on $J$.
REINFORCE -

. Initialize $\theta$ arbitrarily.

. for each episode to

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 $G_t = Q^{\pi_\theta}(S_t, a_t)$

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

$$\theta = \theta + \alpha \, \nabla J(\theta)$$

end

The modification with baseline:

$$\theta_{t+1} \overset{a}{=} \theta_t + a\left(G_t - b(s_t)\right)\frac{\nabla\pi\left(A_t \mid s_t; \theta\right)}{\pi\left(A_t \mid s; \theta\right)}$$

A natural choice of the baseline is an estimate of the state value $\hat{v}(s_t, \omega)$, where $\omega \in \mathbb{R}^d$ is a wt vector learned by a method to be devised!

Softmax policies:) $\pi_\theta(a \mid s) = \dfrac{\exp(\theta_{s,a})}{\sum\limits_{a'} \exp(\theta_{s,a'})}$ $\qquad \textcircled{\scriptsize $\pi$} = \mathbb{R}^{|S||A|}$

2) Linear softmax policies:
$$\pi_\theta(a \mid s) = \frac{\exp(\theta \cdot \phi_{s,a})}{\sum\limits_{a'} \exp(\theta \cdot \phi_{s,a'})}$$

For each $s,a$ we have a feature map $\phi_{s,a} \in \mathbb{R}^d$; then $\theta \in \mathbb{R}^d$.

3) Neural softmax policies:
$$\pi_\theta(a \mid s) = \frac{\exp(f_\theta(s,a))}{\sum\limits_{a'} \exp(f_\theta(s,a'))}$$

$f_\theta(s,a)$ parametrized by a neural network with $\theta \in \mathbb{R}^d$

Recall: in policy gradient, given a trajectory $\tau$, we defined

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

↳ discounted total reward

and

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

and we were optimizing $V^{\pi_\theta}(\mu)$ (i.e.) maximizing expected discounted total reward.

Using softmax policies,

$$\frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta_{s,a}} = \frac{1}{1-\gamma} d^{\pi_\theta}(s) \, \pi_\theta(a|s) A^{\pi_\theta}(s,a)$$

↑ check: easy calculation

- However the gradients with respect to $\theta_{s,a}$ may be small even if $A(s,a)$ is large because $\pi_\theta(s|a)$ is small;

- To prevent probabilities from getting too small

add a regularizer and optimize

$$L_\lambda(\theta) := V^{\pi_\theta}(\mu) - \lambda \underset{s \sim \text{Unif}_S}{\mathbb{E}} \left[ KL\left( \text{Unif}_A, \pi_\theta(\cdot \mid s) \right) \right]$$

$$\textcolor{red}{KL(p,q) := \underset{x \sim p}{\mathbb{E}} - \log\left( \frac{q(x)}{p(x)} \right)}$$

Policy gradient for $L_\lambda(\theta)$:

$$\theta^{t+1} \leftarrow \theta^t + \eta \, \nabla_\theta L_\lambda(\theta^t).$$

THM1   Starting with any initial $\theta^0$ and using

$$\lambda = \frac{\varepsilon(1-\gamma)}{2 \left\| \frac{d_\mu^{\pi^*}}{\mu} \right\|_\infty} \quad \text{and} \quad \eta = 1/\beta\lambda, \quad \beta := \frac{8\gamma}{(1-\gamma)^2} + \frac{2\lambda}{|S|}$$

we have, for all starting distributions $\rho$,

$$\underset{t < T}{\max} \left\{ V^*(s) - V^t(s) \right\} \le \varepsilon, \quad \text{for } T \ge O\left( \frac{|S|^2 |A|}{(1-\gamma)^6 \varepsilon^2} \left\| \frac{d_\rho^{\pi^*}}{\mu} \right\|_\infty^2 \right)$$

# Natural gradients:

**Idea:** $f(x+y) = f(x) + \nabla f(x)^T y + \frac{1}{2} y^T H(f)_x y.$

A stationary pt in the nbd satisfies:

$$\nabla f(x) + H(f) y = 0.$$

$$\boxed{\therefore \; y = H(f)_x^{-1} \nabla f(x).}$$

<span style="color:red">Symm matrix.</span>

- The metric need to measure distances in the nbd is the one induced by the Hessian of f at $x$.

- Amari argued that for policy gradient algorithms, we shouldn't use the standard Euclidean metric;

If the function being optimized is the
loss function of a parameterized distribution
$f(\pi_\theta)$, $f$: loss, $\pi_\theta$ the distribution, Amari
suggested using <span style="color:red">FISCHER INFORMATION MATRIX</span>,
as the choice for the Sym matrix
giving us the botimal form;

- for each $s \in S$, $\pi_\theta(\cdot|s)$ is a distribution.

$$F_\mu^\theta \triangleq \mathop{\mathbb{E}}_{s \sim d_\mu^{\pi_\theta}} \mathop{\mathbb{E}}_{a \sim \pi_\theta(\cdot|s)} \left[ \nabla \log \pi_\theta(a|s) \nabla \log(\pi_\theta(a|s))^\top \right]$$

- Natural policy gradient : NPG

$$\boxed{\theta \leftarrow \theta + \eta \left(F_\mu^\theta\right)^\top \nabla \upsilon^{\pi_\theta}(\mu).}$$

- Motivation for using that? $\uparrow$

For a distribution $\nu$ over $S \times A$, and $\omega \in \mathbb{R}^d$ ($d = $ # parameters), define the compatible function approximator error $L_\nu$ and minimal compatible function approximator error $L_\nu^*$ as:

$$L_\nu(\omega; \theta) = \mathbb{E}_{s,a \sim \nu}\left[\left(A^{\pi_\theta}(s,a) - \omega^T \nabla_\theta \log \pi_\theta(a|s)\right)^2\right]$$

e.g. $\nu(s,a) = d_\rho^{\pi_\theta} \pi_\theta(a|s)$
\underbrace{\phantom{}}

$$\left(A^{\pi_\theta}(s,a) - \omega \nabla_\theta \log \pi_\theta(a|s)\right)^2$$
\underbrace{}_{T}

$$= \left(A^{\pi_\theta}(s,a) - \omega^T \nabla_\theta \log \pi_\theta(a|s)\right)\left(A^{\pi_\theta}(s,a) - \nabla_\theta \log \pi_\theta(a|s)^T \cdot \omega\right)$$

$$\nabla_\omega : -2 \nabla_\theta \log\left(\pi_\theta(a|s)\right) \cdot A^{\pi_\theta}(s,a)$$

$$+ 2 \nabla_\theta \log \pi_\theta(a|s) \cdot \left(\nabla_\theta \log \pi_\theta(a|s)\right)^T \omega$$

For a stationary point:

$$\frac{\nabla_{\omega} :}{\wedge}$$

$$\mathbb{E}_{s \sim d^{\pi_\theta}} \quad \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} \quad - \overset{\pi_\theta}{A}(s,a) \, \nabla_\theta \log \pi_\theta(a|s) \quad +$$

$$\nabla_\theta \log \pi_\theta(a|s) \, \nabla_\theta \log \pi_\theta(a|s)^T \, \omega$$

.

$$= 0 \quad \text{iff} \cdot \qquad \qquad \leftarrow \color{red}{\nabla V^{\widetilde{\pi_\theta}}(\mu)}$$

$$\color{red}{\text{～～～～～～～}}$$

$$\mathbb{E}_{s \sim d^{\pi_\theta}} \quad \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} \quad A^{\pi_\theta}(s,a) \, \nabla_\theta \log \pi_\theta(a|s) \quad =$$

$$\left[ \mathbb{E}_{s \sim d^{\pi_\theta}} \; \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \; \nabla_\theta \log \pi_\theta(a|s) \, \nabla_\theta \log \pi_\theta(a|s)^T \right] \overline{\omega}$$

$$\boxed{\therefore \; \omega = \left( \mathcal{J}_\theta^\mu \right)^T \nabla V^{\widetilde{\pi_\theta}}(\mu).}$$

Letting $\theta \in \mathbb{R}^{S \times A}$ and using softmax policy and NPG

$$\left( \mathcal{J}_\theta^\mu \right)^+ \nabla V^{\pi_\theta}(\mu) = \frac{1}{1-\gamma} A^{\pi_\theta}(s,a) - c(s).$$

$$\color{green}{\uparrow}$$
$$\color{green}{\text{fn of state}}$$

# A simple update

$$\pi^{(t+1)}(a|s) = \pi^t(a|s) \frac{\exp\left(\eta A^t(s,a) / (1-r)\right)}{Z_t(s)}$$

$$= \sum_{a \in A} \pi^t(a|s) \exp\left(\eta \frac{A^t(s,a)}{(1-r)}\right)$$

Thm: For softmax policy class, perform the update $\theta \Leftarrow \theta + \eta (F_\mu^\theta)^\dagger \nabla V^{\pi_\theta}(\mu)$

having $\mu \in \Delta(S)$

$$V^\tau(\mu) \geq V^*(\mu) - \frac{\log(A)}{\eta T} - \frac{1}{(1-r)^2 T}.$$

Set $\eta = (1-r)^2 \log|A|$, then we get to within $\varepsilon$ in $\frac{2}{(1-r)^2 \varepsilon}$ steps

- Recall $\boxed{\text{Eq 1}}$ $\quad \frac{1}{1-\gamma} \sum_s d^\pi(s) \sum_a \frac{\partial \bar{\pi}_\theta(s,a)}{\partial \theta} \cdot \hat{Q}^{\pi_\theta}(s,a)$ $\cdot$

$$ \| $$

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

$$ \text{or} \quad \nabla \log \bar{\pi}_\theta(s,a) $$

Suppose we have $f_\omega : S \times A \to \mathbb{R}$ is an approximator to $\hat{Q}^\pi$, with parameter $\omega$.

<u>Natural</u>: learn $f_\omega$, by following $\pi$ & updating $\omega$ by $\quad \Delta \omega_t \propto \frac{\partial}{\partial \omega} \left[ \hat{Q}^\pi(s_t, a_t) - f_\omega(s_t, a_t) \right]^2$

$$ \propto \left[ \overset{\wedge\pi}{Q}(s_t, a_t) - f_\omega(s_t, a_t) \right] \frac{\partial f_\omega(s_t, a_t)}{\partial \omega} $$

$\qquad\qquad \uparrow$
Some unbiased estimator

So when we converge to a local minima, $\qquad$ **

$$ \sum_s d^\pi(s) \sum_a \pi(s,a) \left[ \hat{Q}^\pi(s,a) - f_\omega(s,a) \right] \frac{\partial f_\omega(s,a)}{\partial \omega} = 0. $$

Suppose we find $f_\omega(s,a)$ s.t: $\qquad\qquad\qquad \Updownarrow$

$$ \nabla V^{\pi_\theta}(\mu) = \frac{1}{1-\gamma} \sum_s d^\pi(s) \sum_a \frac{\partial \pi_\theta(s,a)}{\partial \theta} f_\omega(s,a) \quad * $$

Its clear that choosing $f_w$ s.t

$$\frac{\partial f_w(s,a)}{\partial w} = \frac{1}{\pi(s,a)} \frac{\partial \pi(a,s)}{\partial \theta}, \quad \text{we can satisfy}$$

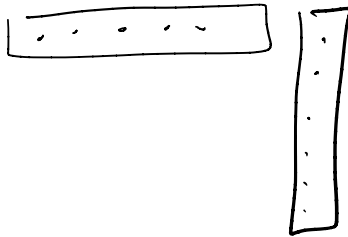<span style="color:red">*</span> <span style="color:green">assuming</span> ** ;

Proof: From ** and $\frac{\partial f_w(s,a)}{\partial w} = \frac{1}{\pi(s,a)} \frac{\partial \pi(a,s)}{\partial \theta}$

we get

$$\sum_s d^\pi(s) \sum_a \frac{\partial \pi(s,a)}{\partial \theta} \left[ Q^\pi(s,a) - f_w(s,a) \right] = 0$$

Eq ②



error is $\perp$ gradient of policy ;

$$\frac{Eq(1) - Eq(2)}{}$$

$$= \sum_s d^\pi(s) \sum_a \frac{\partial \pi(s,a)}{\partial \theta} \left[ Q^\pi(s,a) - Q^\pi(s,a) + f_w(s,a) \right]$$

$$= \sum_s d^\pi(s) \sum_a \frac{\partial \pi(s,a)}{\partial \theta} f_w(s,a)$$

· The gradient has the <span style="color:red">form as in the red expression.</span>

So a good approximator is

$$f_\omega(s,a) = \omega^\top \nabla_\theta \ln \pi_\theta(s,a)$$

then $\dfrac{\partial f_\omega(s,a)}{\partial \omega} \doteq \dfrac{1}{\pi_\theta(s,a)} \nabla_\theta \pi_\theta(s,a)$

· Is the motivation for compatible function approximator definition :-

**Thm!** If the Q value function approximator is compatible to the policy ($\pi$) $\nabla_\omega f_\omega(s,a) = \nabla_\theta \log \pi_\theta(s,a)$

and the parameters $\omega$ minimize

$$\mathbb{E}_{\pi_\theta}\left[ \left( Q^{\pi_\theta}(s,a) - f_\omega(s,a) \right)^2 \right] \quad \text{then}$$

policy gradient is exact

$$\nabla_\theta V^\circ(\mu) = \mathbb{E}_{\pi_\theta}\left[ \nabla_\theta \log \pi_\theta(s,a) \, f_\omega(s,a) \right]$$

- Recall the general value function update

$$v(S_t) \leftarrow v(S_t) \;+\alpha \left[ G_t - V(S_t) \right]$$

<u>Comes from</u>: Minimizing

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

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

- Do a stochastic gradient update at time $t$ since we don't know $v^\pi(s)$;

Get an unbiased estimate — sample $s$ & use $G_t$ instead of $v^\pi(S_t)$

$$\therefore \quad \boxed{v \leftarrow v + \alpha \left( G_t - V(S_t) \right) \frac{\partial v(S_t)}{\partial v}}$$

If $w$ parametrizes values of states,

$$v_w : S \to \mathbb{R},$$

<u>function approximator</u>.

get update: $\boxed{w \leftarrow w + \alpha \left( G_t - v_w(s_t) \right) \dfrac{\partial v_w(s_t)}{\partial w}}$

- ## TD Learning:

Bootstrapping method. bases its update on earlier estimate.

An evaluation alg, estimates $v^{\pi}$

TD update: $v(s) \leftarrow v(s) + \alpha \left( r + \gamma v(s') - v(s) \right)$

(ie) $v(S_t) \leftarrow v(S_t) + \alpha \left[ R_t + \gamma V(S_{t+1}) - V(S_t) \right]$

- Instead of using $G_t$ use $R_t + \gamma V(S_{t+1})$

$$\delta_t = R_t + \gamma V(S_{t+1}) - V(S_t). \quad \text{TD error}$$

- Both TD & MC use estimates: $V(S_{t+1})$ in TD $G_t$ in MC

. Note if we used $v^{\pi}$, $\mathbb{E}\left[\delta_t \mid S_t = s\right] = 0$ .

TD update with function approximator:

$$w \leftarrow w + \alpha \left( R_{t+1} + \gamma v_w(s_{t+1}) - v_w(s_t) \right) \frac{\partial v_w(s_t)}{\partial w}$$

$=$

. TD is not really a gradient update algorithm

. Has good convergence properties; $\rightarrow$ <span style="color:teal">SHOWN TO BE CONVERGENT TO $V_{\pi}$ on an average.</span>

. Converges when using linear function approximation; say $v_w(s) = w^T \underline{\phi(s)}$

features extracted from $s$;

— Comparison with Monte Carlo estimation of $\mathbb{P}_r\left[s, a, s'\right]$ & $R[s, a]$ ?

<span style="color:red">$\updownarrow$</span>

<span style="color:red">TD converges to the same as $v^{\hat{\pi}}$ if every state is seen at least once;</span>

- TD — only $|s|$ values;
- Dynamic / Tabular $|s|^2 |A|$ values;

- SARSA — TD for control:

$$q(s,a) \leftarrow q(s,a) + \alpha \left[ r + \gamma q(s',a') - q(s,a) \right]$$

Stochastic update;

Using function approximator:

$$w \leftarrow w + \alpha \left( r + \gamma q_w(s',a') - q_w(s,a) \right) \frac{\partial q_w}{\partial w} \Big|_{s,a}$$

Given:
- Initialize $q(s,a)$;

for each episode do ($\varepsilon$ greedy / softmax)
 $\quad \pi \leftarrow$ Policy $(q)$
 $\quad s \sim d_0$
 $\quad$ Choose a
 for each time step until $s$ is absorbing do
 $\quad$ Take action $a$, get $r$ & get $s'$
 $\quad$ choose $a'$ from $s'$ using $\pi$
 $\quad q(s,a) \leftarrow q(s,a) + \alpha \left[ r + q(s',a') - q(s,a) \right]$
 $\quad s \leftarrow s'$; $a \in a'$

- Similarly SARSA with function approximator; REINFORCE with baseline for estimating $\pi_\theta \simeq \pi_*$.

Input: A differentiable policy parametrization $\pi(a|s,\theta)$ & state value parametrization $\hat{v}(s,w)$

Parameters: $\alpha^\theta > 0$, $\alpha^w > 0$

Initialize: $\theta$ & $w$

Loop (for each episode):
$\quad S_0, A_0, R_1, \ldots, S_{T-1} A_{T-1} R_T$ following $\pi(\cdot|\cdot;\theta)$
$\quad$ for each step $t = 0, \ldots, T-1$
$$G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$$
$$\delta \leftarrow G - \hat{v}(S_t, w)$$
$$w \leftarrow w + \alpha^w \delta \nabla_w \hat{v}(S_t, w)$$
$$\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla \ln \pi(A_t|S_t,\theta).$$
$\quad$ end.

. Q learning : Off policy TD- control.

○ gets an estimate of $q^*$ ;

$$q(s,a) \leftarrow q(s,a) + \alpha \left[ r + \gamma \max_{a'} q(s',a') - q(s,a) \right]$$

Off policy - update done to get to $q^*$
and really independent of $\pi$ used to generate
the data;

○ Algorith as before;

Next TD($\lambda$). :

Combine the goodness of MC & TD(o)

updates after the
entire episode is run;

immediate update

# n-step bootstrapping.

- Has the benefits of MC & one-step:
  Actions cause updates instantly, in TD(0),
  but variance is large because of this;

- MC- variance is less ∵ updates are done
  at the end.

- In MC updates are performed for each state
  based on the entire sequence of rewards
  from that state to the end of an episode.

- Natural idea: Perform an update based on
  an intermediate number of rewards.

Ex: A two step update, would look at first
    two rewards and estimated value of the
    state two steps later.

- $$G_{t,t+2} := R_{t+1} + \gamma R_{t+2} + \gamma^2 \underbrace{V_{t+1}(S_{t+2})}$$

<span style="color:red">the estimate for $G_{t+2}$ takes place of</span>

<span style="color:green">$\gamma^2 R_{t+3} + \gamma^4 R_{t+4} + \cdots + \gamma^{T-t-1} R_T$ ;</span>

- target for $n$-step update is the $n$-step return

$$G_{t:t+n} := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \cdots + \gamma^{n-1} R_{t+n-1}$$

$$+ \gamma^n V_{t+n-1}(S_{t+n}).$$

$$= ( \text{if } t+n \geq T, \text{ missing terms are treated as zero})$$

- Involves future rewards; to use it we need to $R_{t+n}$ and need to compute $V_{t+n-1}$

- $$\boxed{V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha \left[ G_{t:t+n} - V_{t+n-1}(S_t) \right] \\ 0 \leq t \leq T}$$

$V_{t+n}(s) = V_{t+n-1}(s) \quad \forall \quad s \neq S_t$,

• No changes during first $n-1$ steps of an episode

<span style="color:red">TD[n].</span>

<span style="color:blue">Input: $\pi$</span>
<span style="color:blue">Parameters: $\alpha, n$</span>

Initialize $V(s)$ arbitrarily $\forall \ s \in S$.
Loop for each episode:
Initialize & store $S_0 \neq$ terminal
$T \leftarrow \infty$
Loop for $t = 0, 1, \ldots$
    If $t < T$
        Take action as per $\pi(\cdot | S_t)$
        Observe & store reward as $R_{t+1}$, next state $S_{t+1}$
        If $S_{t+1}$ is terminal $T \leftarrow t+1$
    $\tau \leftarrow t - n + 1$
    If $\tau \geq 0$, $G \leftarrow \displaystyle\sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
    If $\tau + n < T$, $G \leftarrow G + \gamma^n V(S_{\tau+n})$
    $V(S_\tau) \leftarrow V(S_\tau) + \alpha[G - V(S_\tau)]$
until $\tau = T-1$

$$\max_{s} \left[ \mathbb{E}_{\pi} \left[ G_{t:t+n} \mid S_t = s \right] - V_{\pi}(s) \right] \leq$$

$$\gamma^n \max_{s} \left| V_{t+n-1}(s) - v_{\pi}(s) \right|$$

$\forall n \geq 1$