- $k$ : index for behavior policy
- $e$ = index for evaluation policy (fixed)
- $N$ : total number of policies
- $K$ : Number of behavior policies
- $H_k$ : Number of trajectories, $\tau$, in policy $\pi_k$ : Horizon
- $\tau_t = (s_t, a_t, r_t, s_{t+1})$, $s$ = state, $a$ = action, $r$ = reward ; Trajectory
- $\pi_{\theta_k}$ : $[\pi_i$ for $i$ in range $([0 \cdots b, b+1 \cdots N])]$
- $\pi_E$ : dict $\{ \pi_i : \pi_{b \neq i} , \text{ for } i \text{ in range } (N) \}$

- $\rho_k = \prod\limits_{t=1}^{H_k} \dfrac{\pi_e(a_t^k | s_t^k)}{\pi_k(a_t^k | s_t^k)}$ , Importance Sampling ratio

- $\xi_k : [*\rho_k]$ ; i.e $\rho_k$ unwrapped where size $\xi_k = H_k$ and $\xi_k^i = \dfrac{\pi_e(a_i^k | s_i^k)}{\pi_k(a_i^k | s_i^k)}$

- $R_k^i$ : Return (total reward) of $\tau_i$ in policy $k$.

- $\sigma(\pi_e, \pi_k) = \sum\limits_{i=1}^{H_k} R_k^i \times \rho_k^i$

- $X_1 = \dfrac{1}{K \cdot N} \sum\limits_{k=1}^{K} \sigma(\pi_e, \pi_k)$

$$\cdot \; X_2 = \frac{1}{K \cdot N} \sum_{k=1}^{K} \sum_{i=1}^{N} R_k^i$$

$$\cdot \; V(\pi_e) = \frac{1}{N} \sum_{i=1}^{N} R_e^i \quad \in \mathbb{R}^N \text{ (vector of values of evaluation policies for N policies)}$$

$\cdot \; R_e^i$ : Total reward of the trajectory $i$ of evaluation policy $e$.

$$\cdot \; Error = \left( \underbrace{[X_1 \; X_2 \; 1] \begin{bmatrix} c \\ d \\ e \end{bmatrix}}_{\text{estimated value function}} - \underbrace{V(\pi_e)}_{\text{true Value function}} \right)^2$$

---

$$\boxed{\text{Sampling Trajectories} \; 2}$$

Define: $k, e, \text{model}_\pi, H$ :    // model $\Rightarrow$ TD / Monte Carlo / PPO / DQN / etc —

def getAction(state, Q): // get discrete action & associated probability (softmax)

     return $\text{argmax}(Q[\text{state}])$, $\max \left\{ \dfrac{\exp(Q[\text{state}])}{\sum \exp(Q[\text{state}])} \right\}$

R=0 // reward total

for $i$ in range(H):

     $a, \text{prob}_k = \text{getAction}(s_t, \text{model}_{\pi_k})$;

     $-, \text{prob}_e = \text{getAction}(s_t, \text{model}_{\pi_e})$;

     // note that behavior policy samples trajectories, $\tau$)

     $s_{t+1}, r = \pi_k(a_t|s_t)$ // $\pi_k$ // generates next state & reward from current state, $s_t$

     R = R+r

   --- // code

- Assumes IS weights are already calculated (see (1)). This section uses 1-step gradient update using MAML algorithm. Note: this is a single-task problem

- <u>Goal</u>: find set of importance sampling weights, $w$, such that $\rho \to 1$.

This is also equivalent to finding parameters $w_\kappa$ from $\pi_\kappa$ such that

$$\sum_{i=1}^{H_\kappa} \frac{\pi_\theta(a_i^\kappa | s_i^\kappa)}{\pi_\kappa(a_i^\kappa | s_i^\kappa)} \to H_\kappa$$

let weights, $w, = \left[ \frac{1}{\pi_\kappa(a_1^\kappa | s_1^\kappa)}, \frac{1}{\pi_\kappa(a_2^\kappa | s_2^\kappa)}, \cdots, \frac{1}{\pi_\kappa(a_{H_\kappa}^\kappa | s_{H_\kappa}^\kappa)} \right] \in \mathbb{R}^{H_\kappa}$

let feature matrix, $X = \left[ \pi_\theta(a_1^\kappa | s_1^\kappa), \pi_\theta(a_2^\kappa | s_2^\kappa), \cdots, \pi_\theta(a_{H_\kappa}^\kappa | s_{H_\kappa}^\kappa) \right] \in \mathbb{R}^{H_\kappa}$

$\Rightarrow w \cdot X = H_\kappa = [w_1, \cdots, w_{H_\kappa}] \begin{bmatrix} X_1 \\ \vdots \\ X_{H_\kappa} \end{bmatrix} = H_\kappa \quad \Rightarrow (1 \times H_\kappa)(H_\kappa \times 1) = \underline{\underline{1 \times 1}} \\ \in \mathbb{R}^1$

- $X \to$ fixed, $w \to$ parameters to optimize.

Set MAML framework such that $\quad \underline{\underline{MAML(w) \cdot X - H_\kappa < w \cdot X - H_\kappa}}$

MAML $\begin{cases} \text{while not done do:} \\ \qquad \text{Evaluate } \nabla_w \mathcal{L}_X(f(w)) \\ \qquad \text{Compute adapted parameters w/ SGD: } w_i' = w - \alpha \nabla_w \mathcal{L}_X(f(w)) \\ \\ \qquad w \leftarrow w - \beta \nabla_w \sum \mathcal{L}_X(f w_i') \end{cases}$

Example 4

- let $N = 3$

→ data point 1 :

$$\pi_{R} = \{\pi_2, \pi_3\} , \quad \pi_e = \pi_1$$

→ data point 2 :

$$\pi_{K} = \{\pi_1, \pi_3\} , \pi_e = \pi_2$$

→ data point 3 :

$$\pi_{R} = \{\pi_1, \pi_2\} , \pi_e = \pi_3$$

- Sample trajectories based on $\pi_K$ & $\pi_e$ (see (2))
- MAML ($\pi_K$), see (3)
- Compute $\rho_R$
- Compute $\sigma(\pi_e, \pi_R)$
- Compute $V(\pi_e)$
- Compute $X_1, X_2$
- Find MSE.;  $\left( [x_1 \; x_2 \; 1] \begin{bmatrix} c \\ d \\ e \end{bmatrix} - V(\pi_e) \right)^2$ — ✱
- Optimize parameters $c, d,$ and $e$ such that (✱) $= min$