

↑; DR FROM PREVIOUS LECTURE

•) GPI \equiv ESTIMATE OF q_{π} \rightarrow IMPROVEMENT OF π USING q_{π}
($\pi' = \epsilon$ -greedy(q_{π})) q_{π}

↑
WHY SOLVE A MORE GENERAL PROBLEM?



IDEA: ① $\forall \alpha \in S$, PARAMETRIZE THE POLICY: $\pi(a, s) = \pi(a, s, \theta) := \pi(\theta)$

② CHOOSE A PERFORMANCE MEASURE $J(\theta)$

THE LARGER $J(\theta)$, THE BETTER $\pi(\theta)$

③ OPTIMIZE $J(\theta)$: $\theta^* := \arg \max_{\theta} J(\theta)$

① SOFTMAX PARAMETRIZATION \equiv COMPUTE PREFERENCES $h(a, s, \theta)$ OF STATE-ACTION PAIRS (s, a) (FOR INSTANCE, USING A NN), AND DEFINE ACTION PROBABILITIES BY:

$$\pi(a, s, \theta) := \frac{e^{h(a, s, \theta)}}{\sum_{b \in A} e^{h(a, b, \theta)}}$$

② $J(\theta) =$ "EXPECTED TOTAL REWARD GIVEN BY $\pi(\theta)$ "
 $=$ "TOTAL REWARD ALONG TRAJECTORIES τ WEIGHTED BY $\text{PROB}(\tau)$ "

•) $\tau = s_0, a_0, r_1, s_1, a_1, \dots$

•) $\text{PROB}(\tau) = p(s_0) \cdot \pi(a_0, s_0, \theta) \cdot p(r_1, s_1, s_0, a_0) \cdot \pi(a_1, s_1, \theta) \cdot \dots$

STARTING DISTRIBUTION

DISTRIBUTION MODEL

•) "TOTAL REWARD ALONG τ " $\equiv \sum_{t=0}^{\infty} \gamma^t r_{t+1} =$ "DISCOUNTED RETURN" $:= R(\tau)$

•) "EXPECTED TOTAL REWARD" \equiv

$$J(\theta) := \sum_{\tau} \text{PROB}(\tau) \cdot R(\tau) = \mathbb{E}_{\tau \sim \text{PROB}_{\theta}} [R(\tau)]$$

REMARK: IF $d_{s_0}^H(\varphi) =$ "PROBABILITY OF REACHING φ FROM s_0 IN H STEPS",

THEN $d_{s_0}(\varphi) := \sum_{H=0}^{\infty} d_{s_0}^H(\varphi) =$ "AVERAGE NUMBER OF VISITS TO φ "

GIVES A DISTRIBUTION ON \mathcal{S} : $\mu(\varphi) = \frac{d_{s_0}(\varphi)}{\sum_{\varphi} d_{s_0}(\varphi)}$

↖ "ON-POLICY DISTRIBUTION"

↙ \equiv "PROBABILITY OF VISITING & STARTING FROM s_0 "

THEOREM

$J(\theta)$ CAN BE COMPUTED BY LOOKING ONE

STEP FORWARD :

$$J(\theta) = \sum_{s,a} \mu(s) \cdot \pi(a,s,\theta) \cdot r(s,a)$$

"DISTRIBUTIONAL POINT OF VIEW"

EXERCISE: WRITE THIS
WHEN STARTING FROM
A FIXED s_0 .

$$= \sum_s \rho(s) \cdot \left(\sum_a \pi(a,s,\theta) \cdot r_{\pi_\theta}(s,a) \right)$$

③ OPTIMIZE $J(\theta) \rightsquigarrow$ WE NEED DERIVATIVES $\nabla_{\theta} J = \nabla J$

"GRADIENT OF J WRT θ "

THIS LECTURE!