

APPENDIX A
PROOF OF THEOREM 1

Theorem 1 (The Generalized Policy Gradient Theorem). *The derivative of $J_s(\pi_\theta)$ with respect to θ is the expectation of the product of the π_θ -induced trajectory's SOTA probability and the gradient of the log of policy π_θ , i.e., $\nabla_\theta J_s(\pi_\theta) =$*

$\mathbb{E}_{\tau \sim \pi} [\mathbb{P}[R(\tau) \leq T] \nabla_\theta \log \pi_\theta(\tau)]$, where τ is π_θ -induced trajectory, $R(\tau)$ is τ 's total travel time (an RV), and $\pi_\theta(\tau)$ refers to the probability of generating τ by π_θ .

Proof.

$$\begin{aligned}
 \nabla_\theta J_s(\pi_\theta) &= \nabla_\theta \mathbb{P}[t^\pi(o, d) \leq T] \\
 &= \nabla_\theta \int_\tau \mathbb{P}[R(\tau) \leq T] \pi_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \nabla_\theta \pi_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} \pi_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \nabla_\theta \log \pi_\theta(\tau) \pi_\theta(\tau) d\tau \\
 &= \mathbb{E}_{\tau \sim \pi_\theta} [\mathbb{P}[R(\tau) \leq T] \nabla_\theta \log \pi_\theta(\tau)].
 \end{aligned}$$

□

In the proof process, $\pi_\theta(\tau)$ refers to the probability that trajectory τ is generated by π_θ .

APPENDIX B
PROOF OF THEOREM 2

Theorem 2 (Off-Policy Generalized Policy Gradient Theorem). *The gradient of $J_s(\pi_\theta)$ for behavior policy (μ_θ) generated trajectory τ can be expressed as:*

$$\nabla J_s(\pi_\theta) = \mathbb{E}_{\tau \sim \mu_\theta} [\mathbb{P}[R(\tau) \leq T] (\nabla_\theta \log \pi_\theta(\tau)) \rho(\tau)], \quad (14)$$

where $\rho(\tau) = \pi_\theta(\tau) / \mu_\theta(\tau)$ is the importance sampling ratio of π to μ with respect to trajectory τ .

Proof.

$$\begin{aligned}
 \nabla_\theta J_s(\pi_\theta) &= \nabla_\theta \mathbb{P}[t^\pi(o, d) \leq T] \\
 &= \nabla_\theta \int_\tau \mathbb{P}[R(\tau) \leq T] \pi_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \nabla_\theta \pi_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} \frac{\pi_\theta(\tau)}{\mu_\theta(\tau)} \mu_\theta(\tau) d\tau \\
 &= \int_\tau \mathbb{P}[R(\tau) \leq T] \nabla_\theta \log \pi_\theta(\tau) \rho(\tau) \mu_\theta(\tau) d\tau \\
 &= \mathbb{E}_{\tau \sim \mu_\theta} [\mathbb{P}[R(\tau) \leq T] (\nabla_\theta \log \pi_\theta(\tau)) \rho(\tau)].
 \end{aligned}$$

□