

Abstract geometric lines in black on a white background, forming various overlapping polygons and shapes.

# SOFT ACTOR-CRITIC

Off-Policy Maximum Entropy Deep Reinforcement  
Learning with a Stochastic Actor

Haarnoja et al.

Alireza Nobakht

Deep Learning Course • Dr. Samaneh Hosseini

Isfahan University of Technology

# OUTLINE

Introduction

Related work

Soft Actor-Critic

New objective

Policy update rule

Cost functions and gradients

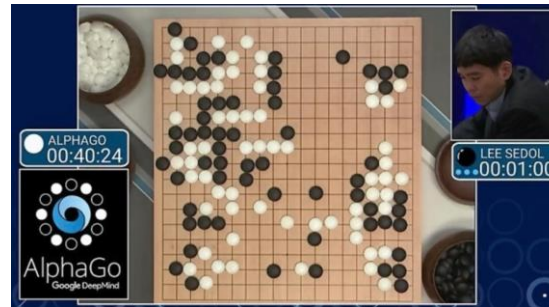
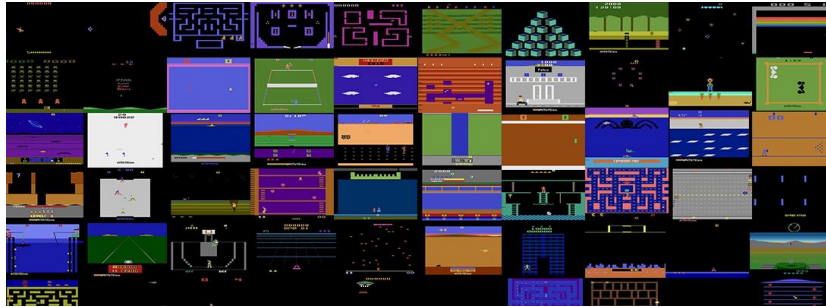
Algorithm

Experiments

# INTRODUCTION

## Model-Free Reinforcement Learning

- In combination with high-capacity function approximators such as neural networks



# INTRODUCTION

## Model-Free Reinforcement Learning Challenges

- Very High Sample Complexity
  - Requires millions of steps
- brittle with respect to their hyperparameters
  - Learning rate, exploration constants and ...
- Continuous state and action spaces

# RELATED WORK

## Off-policy

- E.g. DDPG
- Improve Sample Complexity
- Extremely difficult to stabilize and brittle to hyperparameter settings
  - difficult to extend to complex, high-dimensional tasks
- Continuous state and action spaces

## On-policy

- E.g. PPO
- Improve Stability
- Good for continuous state and action spaces
- Poor sample complexity



# SOFT ACTOR-CRITIC

an actor-critic architecture with separate policy and value function networks

an off-policy formulation that enables reuse of previously collected data for efficiency

entropy maximization to enable stability and exploration

# NEW OBJECTIVE

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} [r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_t))]$$

$$J(\pi) = \sum_{t=0}^{\infty} \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} \left[ \sum_{l=t}^{\infty} \gamma^{l-t} \mathbb{E}_{\mathbf{s}_l \sim p, \mathbf{a}_l \sim \pi} [r(\mathbf{s}_l, \mathbf{a}_l) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_l)) | \mathbf{s}_t, \mathbf{a}_t] \right]$$

$\alpha$ : The temperature parameter

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

# POLICY UPDATE RULE

$$\pi_{\text{new}} = \arg \min_{\pi' \in \Pi} D_{\text{KL}} \left( \pi'(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot))}{Z^{\pi_{\text{old}}}(\mathbf{s}_t)} \right)$$

**Lemma 2** (Soft Policy Improvement). *Let  $\pi_{\text{old}} \in \Pi$  and let  $\pi_{\text{new}}$  be the optimizer of the minimization problem defined in [Equation 4](#). Then  $Q^{\pi_{\text{new}}}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t)$  for all  $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$  with  $|\mathcal{A}| < \infty$ .*



# COST FUNCTIONS & GRADIENTS

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ \frac{1}{2} \left( V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)] \right)^2 \right] \quad \hat{\nabla}_\psi J_V(\psi) = \nabla_\psi V_\psi(\mathbf{s}_t) (V_\psi(\mathbf{s}_t) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t) + \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t))$$


---

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right] \quad \hat{\nabla}_\theta J_Q(\theta) = \nabla_\theta Q_\theta(\mathbf{a}_t, \mathbf{s}_t) \left( Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - r(\mathbf{s}_t, \mathbf{a}_t) - \gamma V_{\bar{\psi}}(\mathbf{s}_{t+1}) \right)$$

$$\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V_{\bar{\psi}}(\mathbf{s}_{t+1})]$$


---

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ \text{D}_{\text{KL}} \left( \pi_\phi(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q_\theta(\mathbf{s}_t, \cdot))}{Z_\theta(\mathbf{s}_t)} \right) \right] \quad \hat{\nabla}_\phi J_\pi(\phi) = \nabla_\phi \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$$

$\mathbf{a}_t = f_\phi(\epsilon_t; \mathbf{s}_t)$  reparameterization trick

$$+ (\nabla_{\mathbf{a}_t} \log \pi_\phi(\mathbf{a}_t | \mathbf{s}_t) - \nabla_{\mathbf{a}_t} Q(\mathbf{s}_t, \mathbf{a}_t)) \nabla_\phi f_\phi(\epsilon_t; \mathbf{s}_t)$$

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} [\log \pi_\phi(f_\phi(\epsilon_t; \mathbf{s}_t) | \mathbf{s}_t) - Q_\theta(\mathbf{s}_t, f_\phi(\epsilon_t; \mathbf{s}_t))]$$

# ALGORITHM

Initialize parameter vectors  $\psi, \bar{\psi}, \theta, \phi$ .

**for** each iteration **do**

**for** each environment step **do**

$$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$$

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$$

**end for**

**for** each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$$

$$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i) \text{ for } i \in \{1, 2\}$$

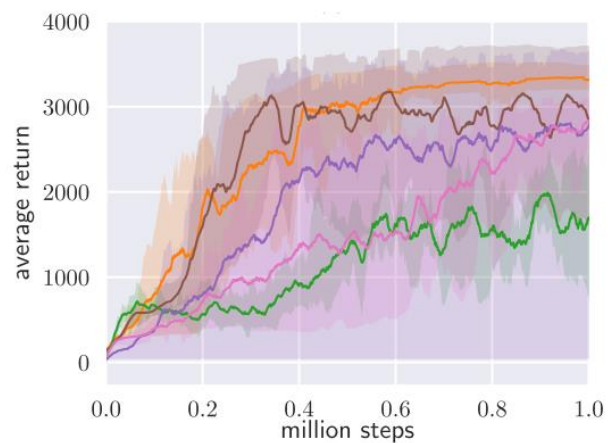
$$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$$

$$\bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}$$

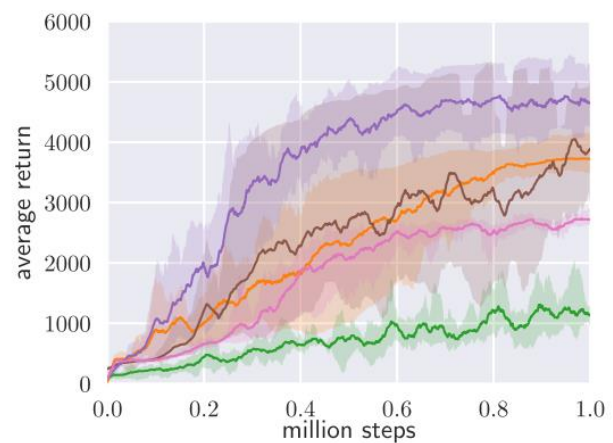
**end for**

**end for**

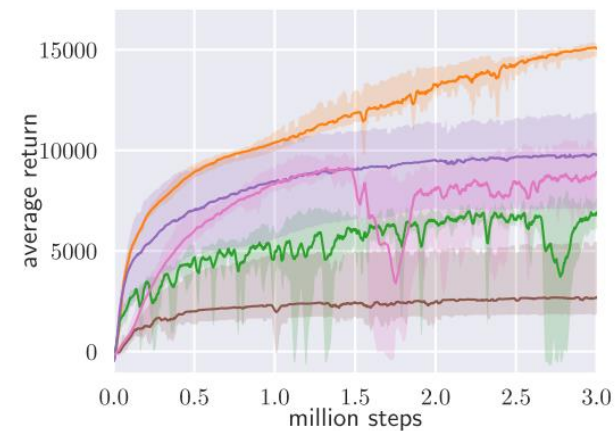
# EXPERIMENTS



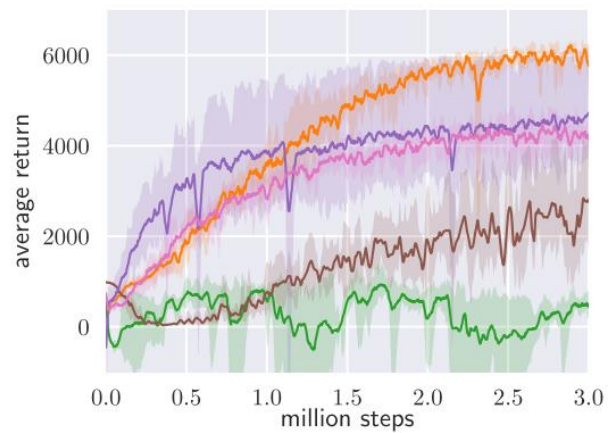
(a) Hopper-v1



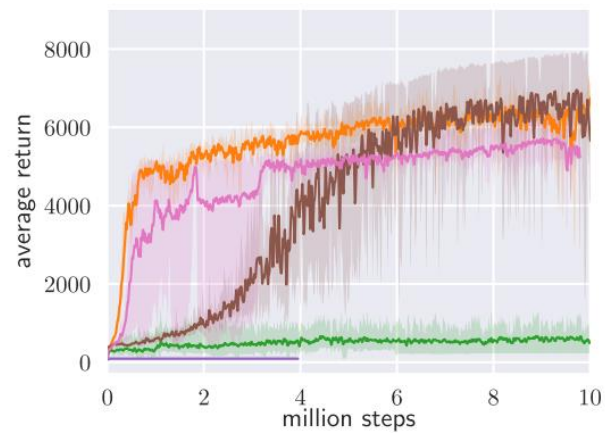
(b) Walker2d-v1



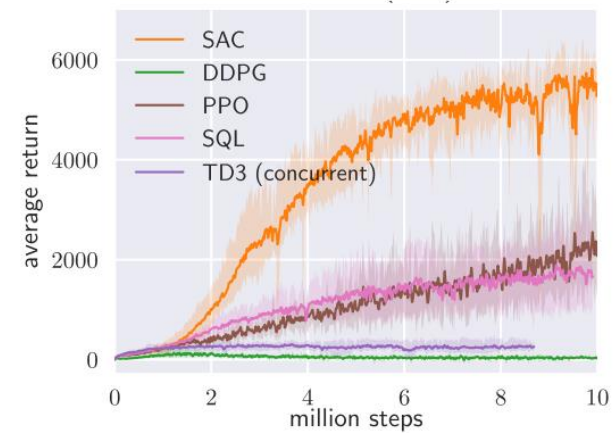
(c) HalfCheetah-v1



(d) Ant-v1

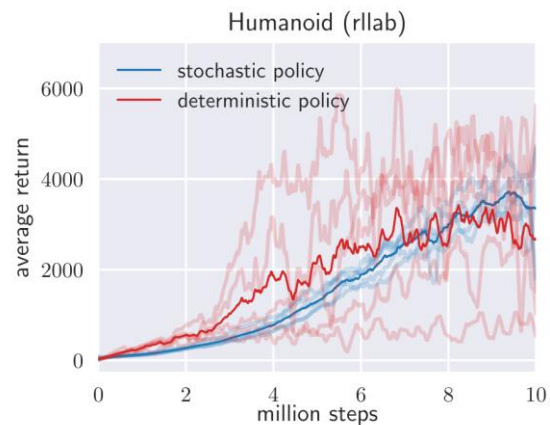


(e) Humanoid-v1

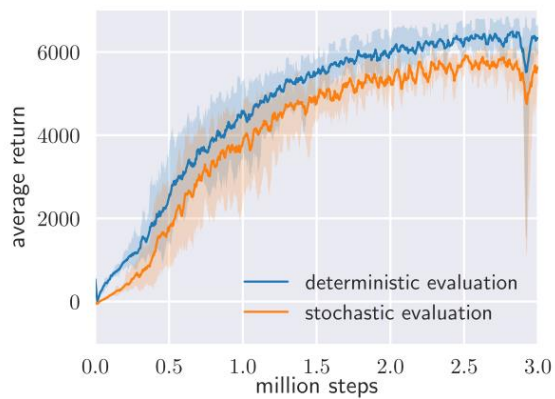


(f) Humanoid (rllab)

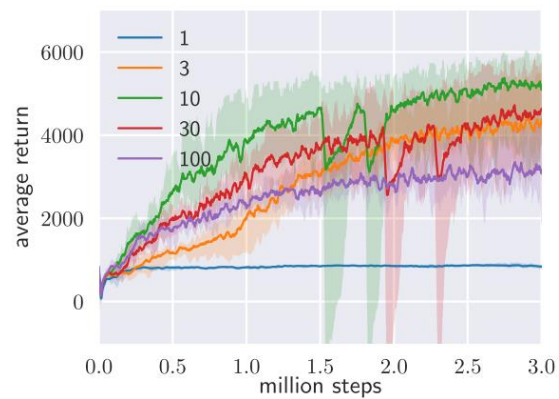
# EXPERIMENTS



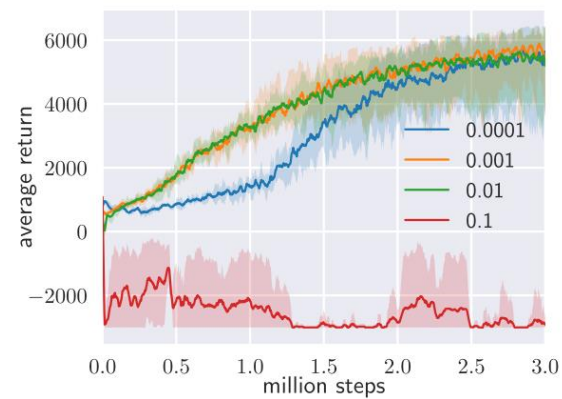
learning a stochastic policy with entropy maximization can drastically stabilize training



(a) Evaluation



(b) Reward Scale



(c) Target Smoothing Coefficient ( $\tau$ )



# THANK YOU

Alireza Nobakht

Deep Learning Course • Dr. Samaneh Hosseini

Isfahan University of Technology