

Map



Model
Based

Map

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy
Off Policy

MLMBTRL
(learn T, R)

Q-learning

L

Model
Based

Map

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy
↓
Off Policy

MLMBTRL
(learn T, R)

Q-learning

Actor



Model
Based

Map

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy

Off Policy

MLMBTRL
(learn T, R)

Q-learning

Actor - Critic



Model
Based

Map

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy
↓
Off Policy

MLMBTRL
(learn T, R)

Q-learning

Actor - Critic

Challenges:

1. Exploration vs Exploitation ↙
2. Credit Assignment
3. Generalization

Model
Based

Map

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy
↓
Off Policy

MLMBTRL
(learn T, R)

Q-learning

Actor - Critic

Challenges:

1. Exploration vs Exploitation ←
2. Credit Assignment
3. Generalization

Is Exploration Important? Montezuma's Revenge

Recognize + track

UCB

Generalization?

Is Exploration Important? Theory

Exploration Bonus

Exploration Bonus

- In General, $R^+(s, a) = R(s, a) + B(s, a)$
- UCB: $B(s, a) = c \sqrt{\frac{\log N(s)}{N(s,a)}}$

Exploration Bonus

Example 1: Learn Pseudocount

Exploration Bonus

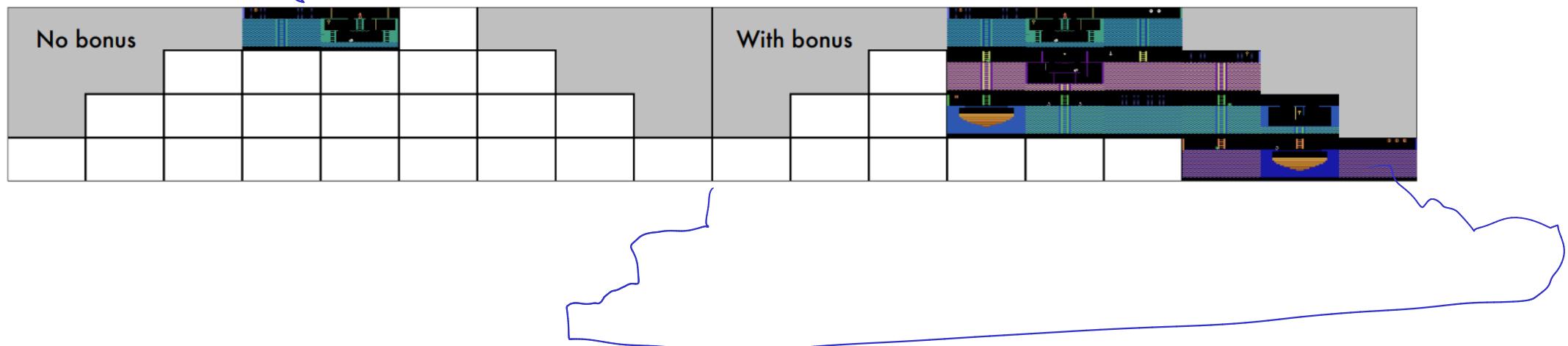
Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation

Exploration Bonus

Example 1: Learn Pseudocount

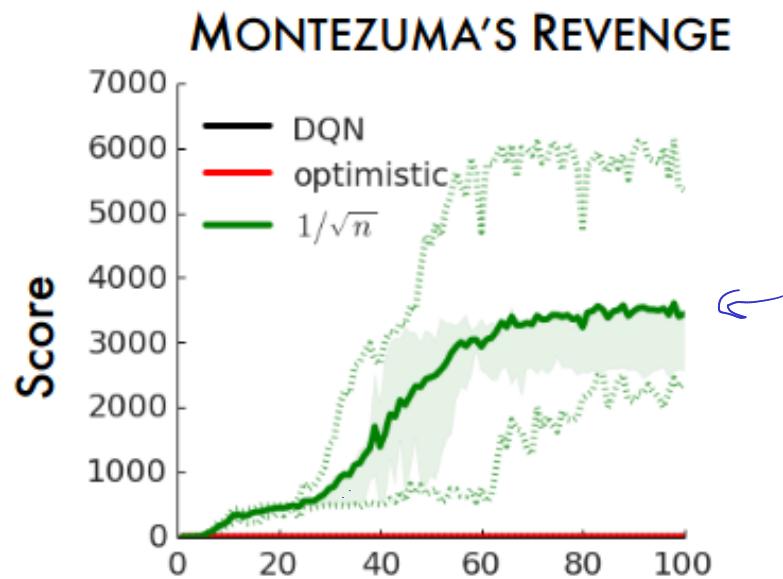
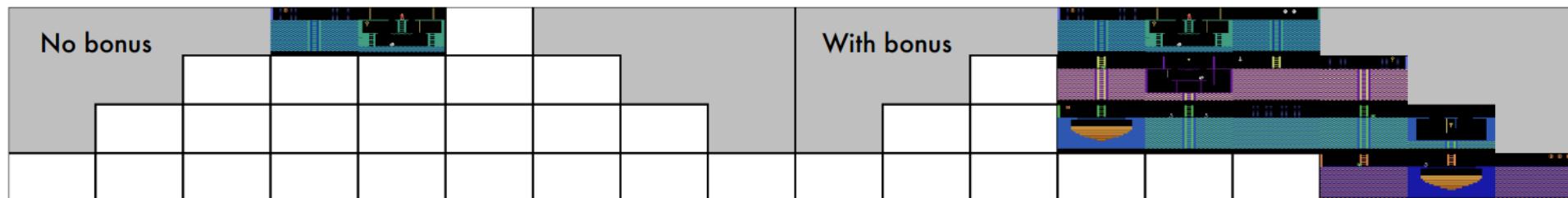
$$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}} \text{ where } \hat{N}(s) \text{ is a learned function approximation}$$



Exploration Bonus

Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation

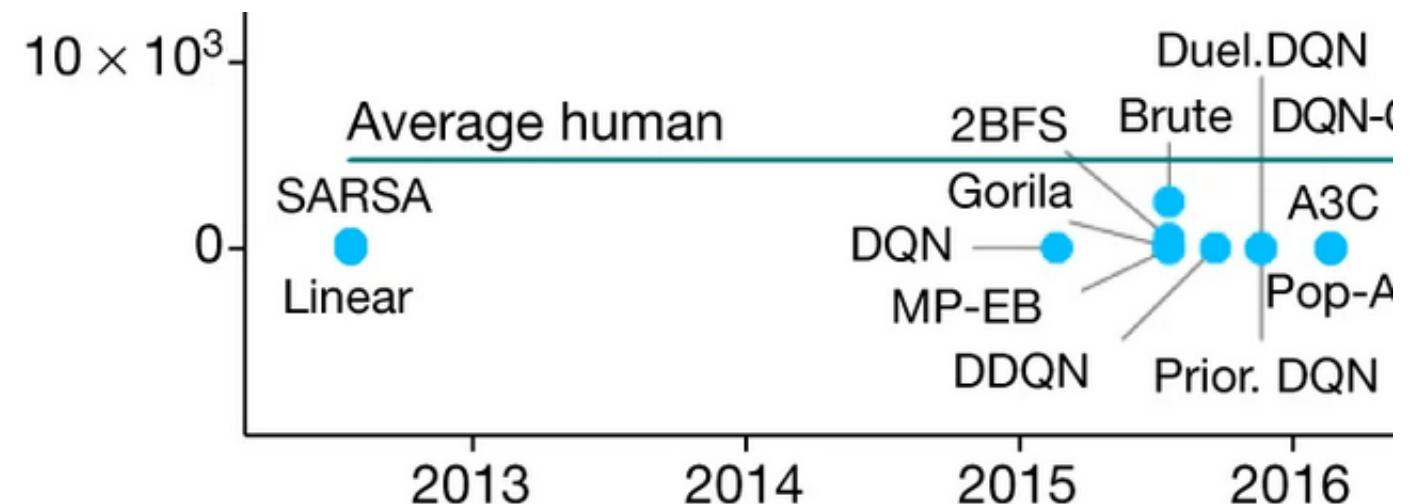
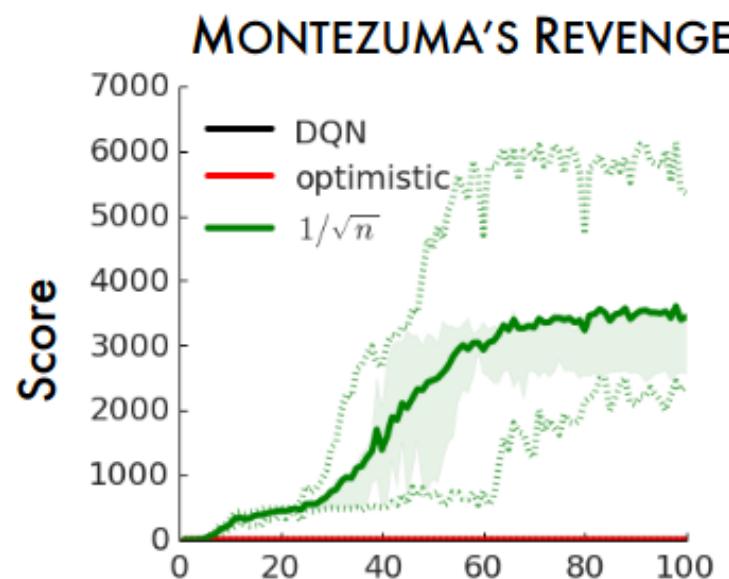
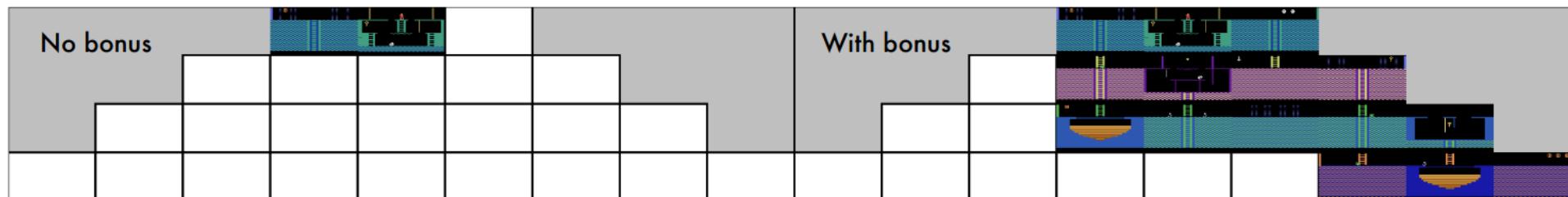


Bellemare, et al. 2016 "Unifying Count-Based Exploration..."

Exploration Bonus

Example 1: Learn Pseudocount

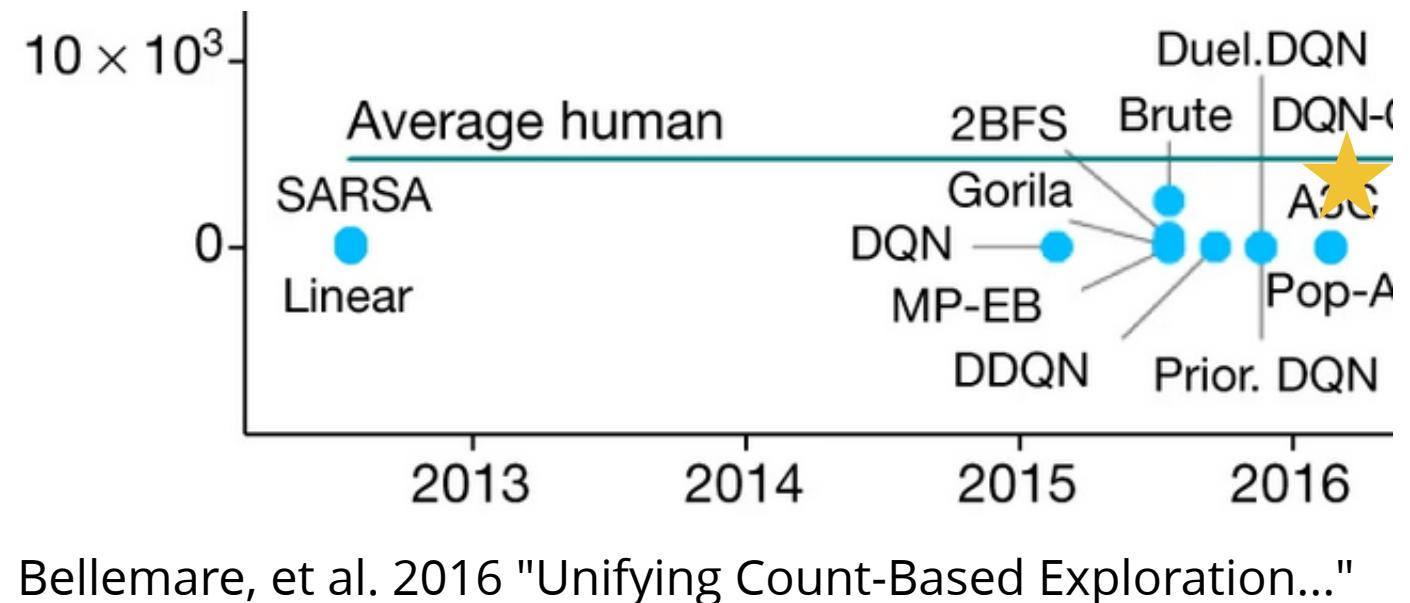
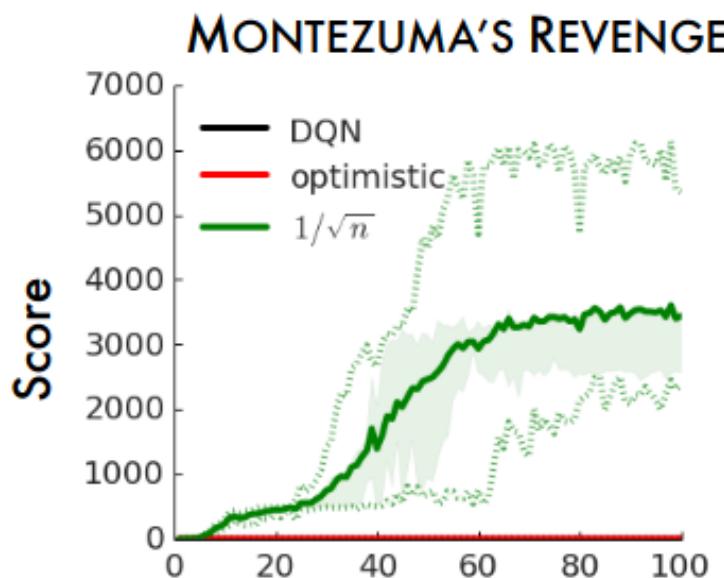
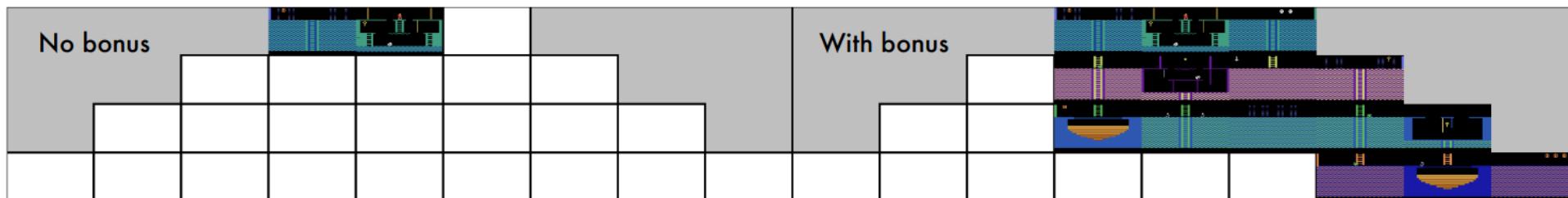
$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation



Exploration Bonus

Example 1: Learn Pseudocount

$B(s, a) \approx \frac{1}{\sqrt{\hat{N}(s)}}$ where $\hat{N}(s)$ is a learned function approximation



Exploration Bonus

^

-

,

Exploration Bonus

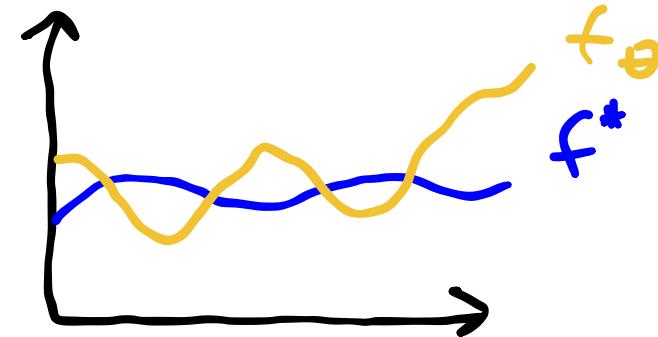
Example 2: Learn a function of the state and action



Exploration Bonus

Example 2: Learn a function of the state and action

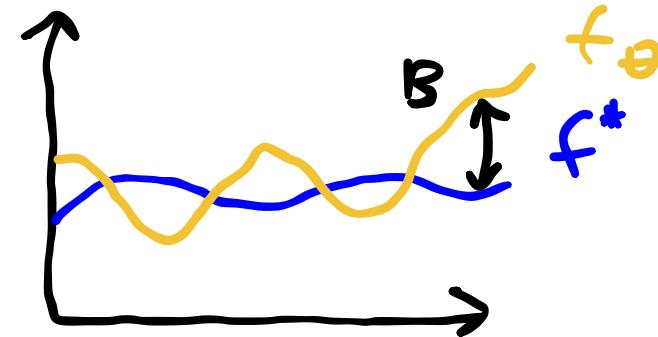
$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$



Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

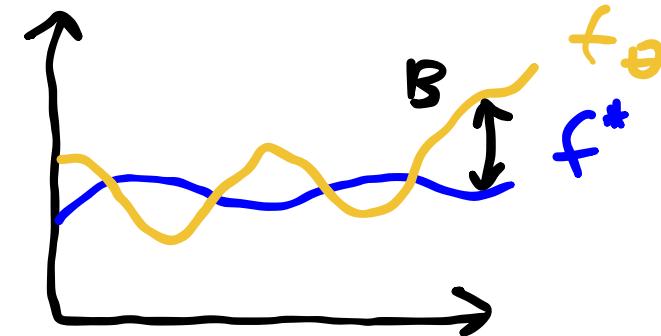


Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

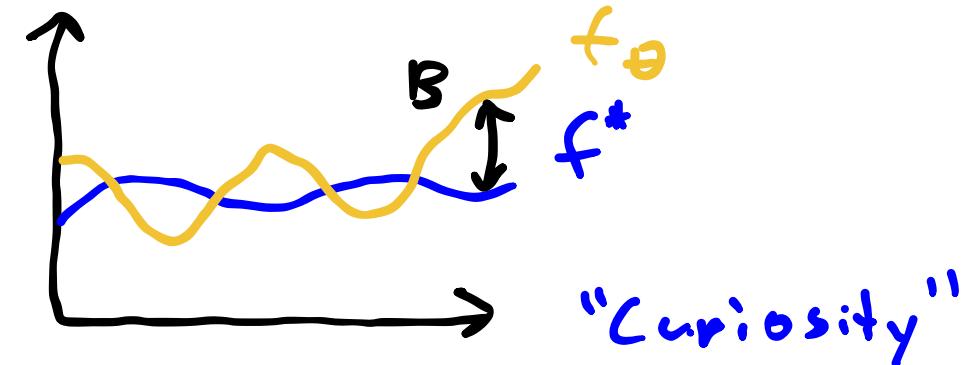


Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?



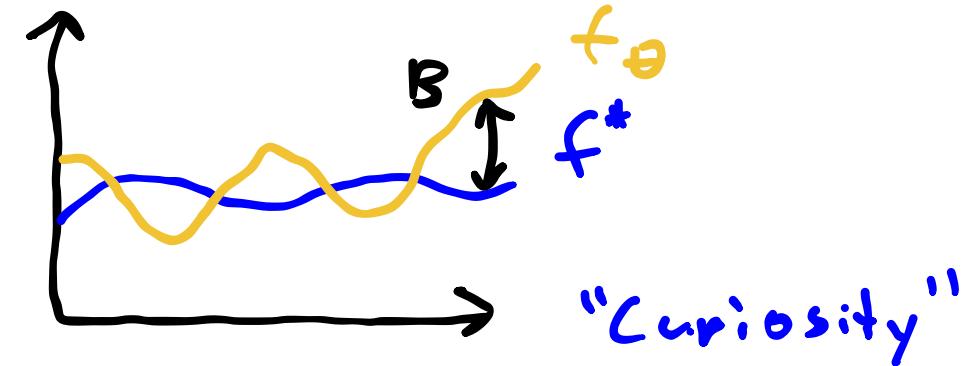
Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)



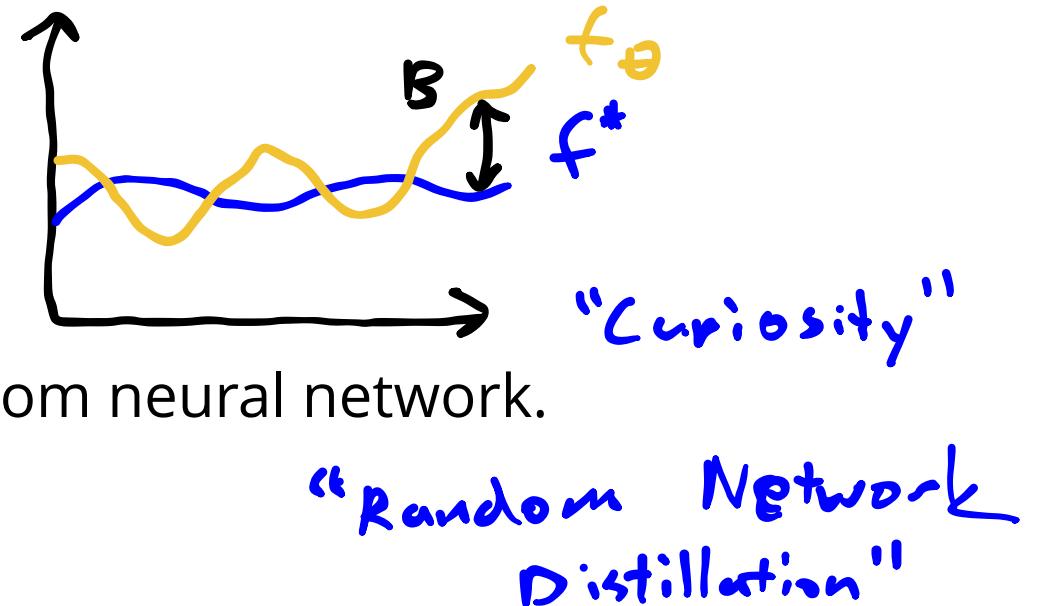
Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)
- $f^*(s, a) = f_\phi(s, a)$ where f_ϕ is a random neural network.



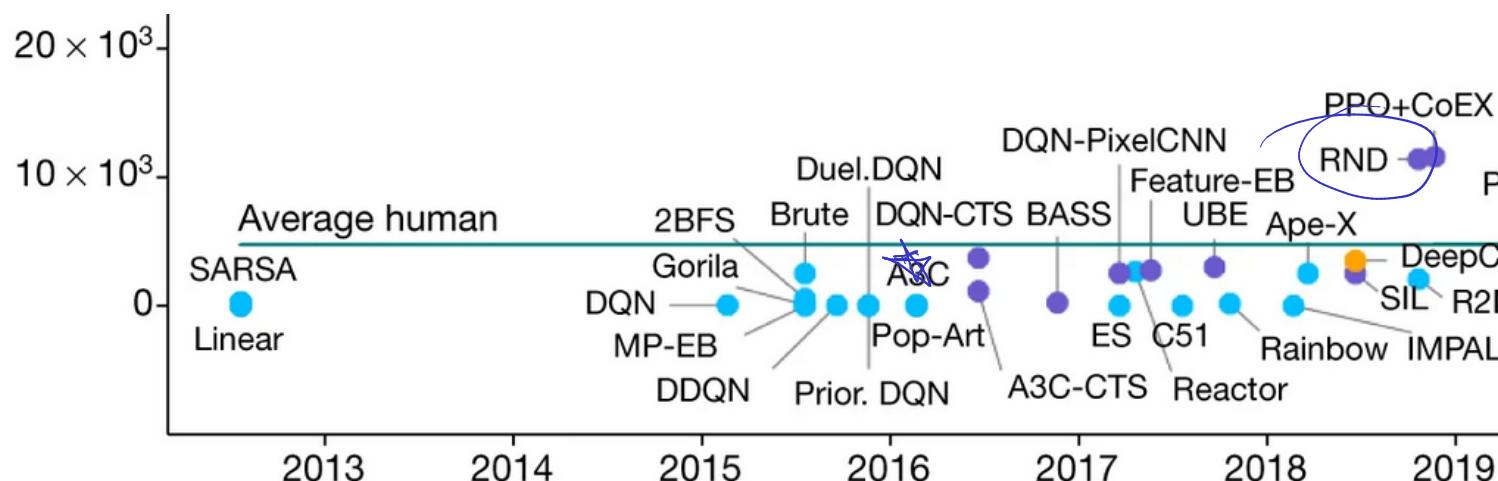
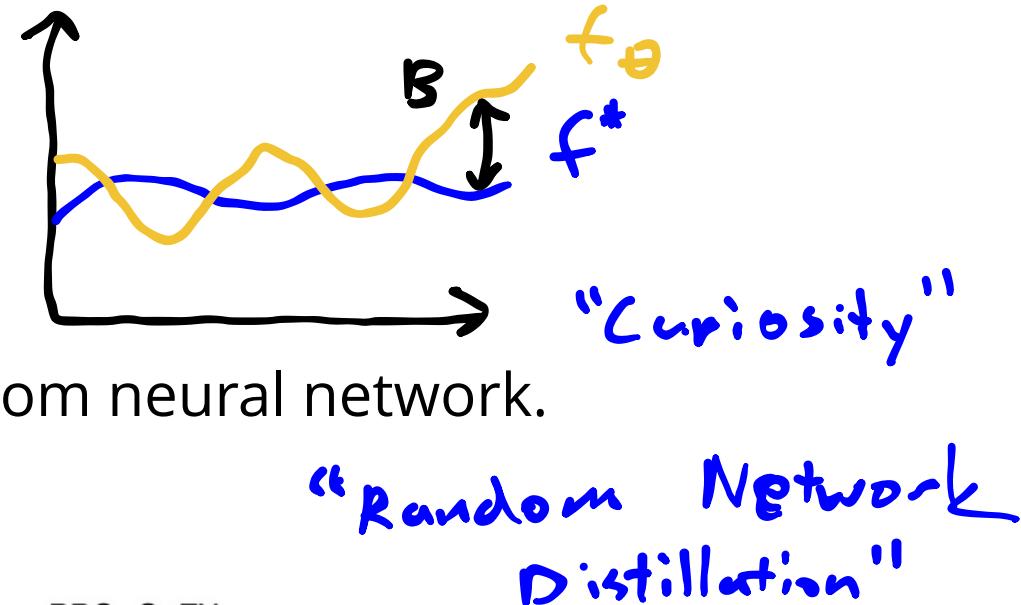
Exploration Bonus

Example 2: Learn a function of the state and action

$$B(s, a) = \|\hat{f}_\theta(s, a) - f^*(s, a)\|^2$$

What should f^* be?

- $f^*(s, a) = s'$ (Next state prediction)
- $f^*(s, a) = f_\phi(s, a)$ where f_ϕ is a random neural network.



Exploration Bonus

Exploration Bonus

Example 3: Thompson Sampling

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q
2. Sample Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q
2. Sample Q
3. Act according to ~~Q~~ sampled Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q  Hard
2. Sample Q
3. Act according to Q

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q  **Hard**
 2. Sample Q
 3. Act according to Q
-
- Bootstrapping with multiple Q networks

Exploration Bonus

Example 3: Thompson Sampling

1. Maintain a distribution over Q  **Hard**
2. Sample Q
3. Act according to Q

- Bootstrapping with multiple Q networks
- Dropout

Exploration Bonus

Exploration Bonus

Example 4: Go-Explore

Exploration Bonus

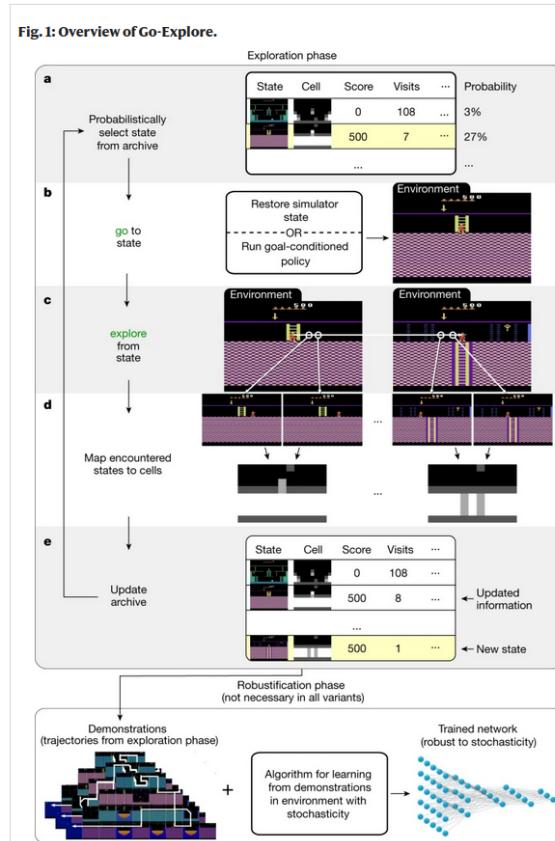
Example 4: Go-Explore

"First return, then explore"

Exploration Bonus

Example 4: Go-Explore

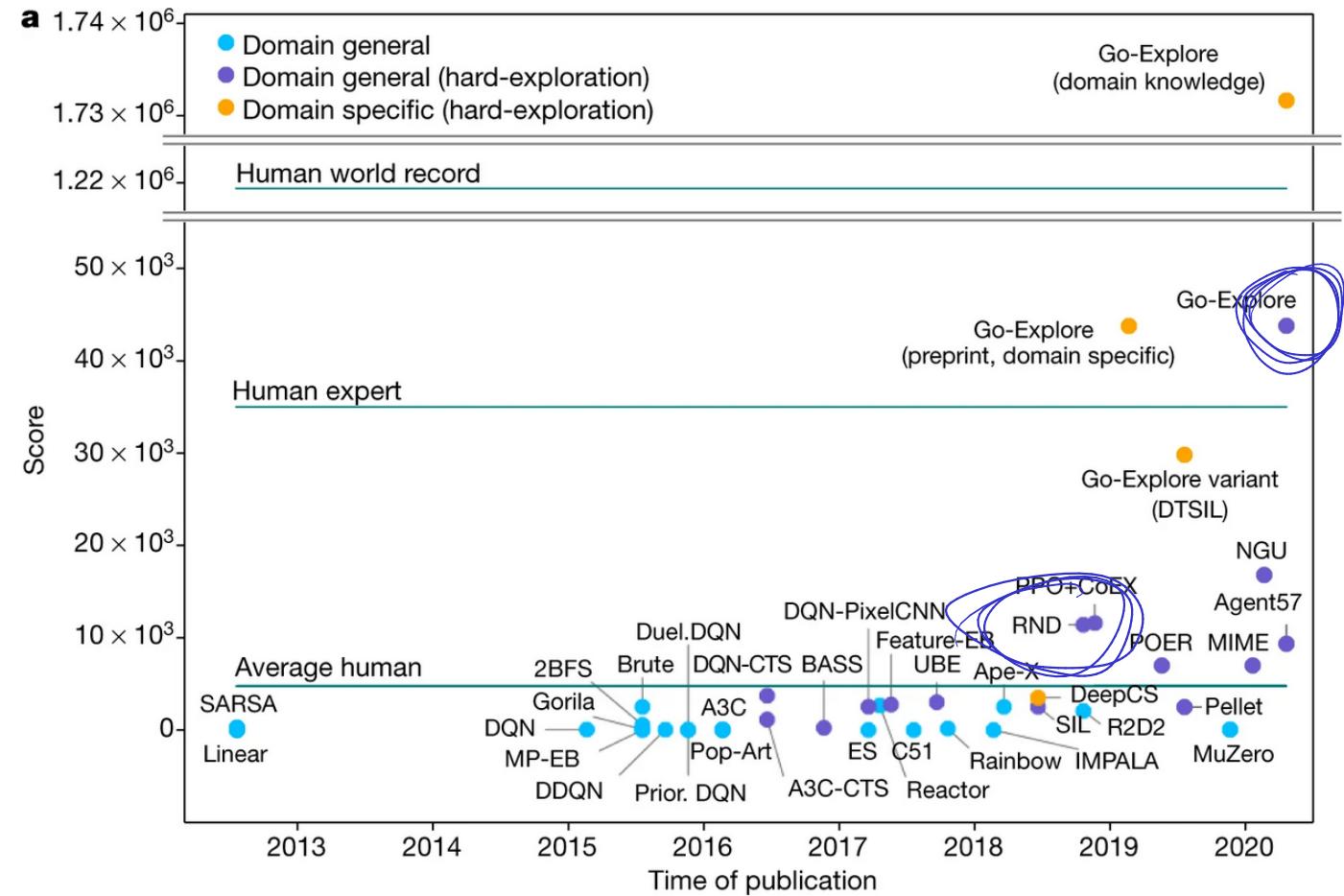
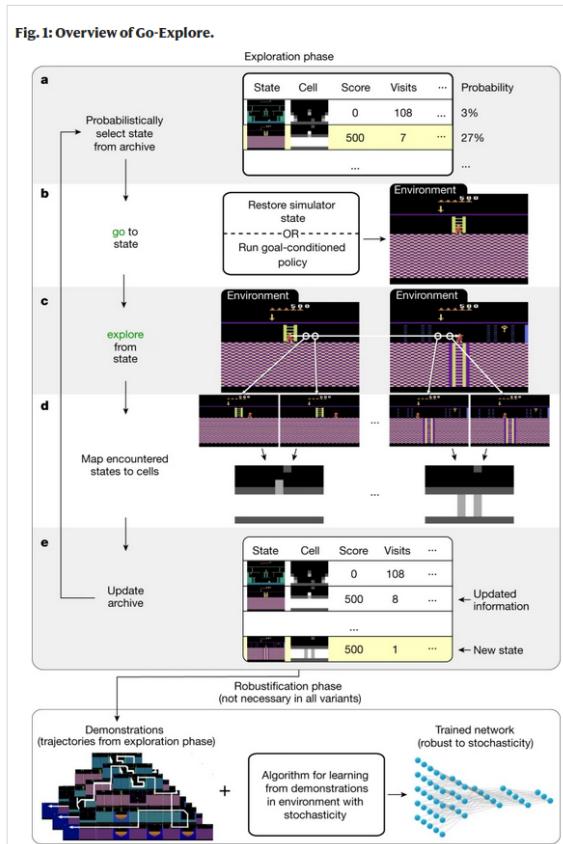
"First return, then explore"



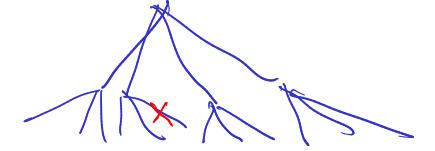
Exploration Bonus

Example 4: Go-Explore

"First return, then explore"

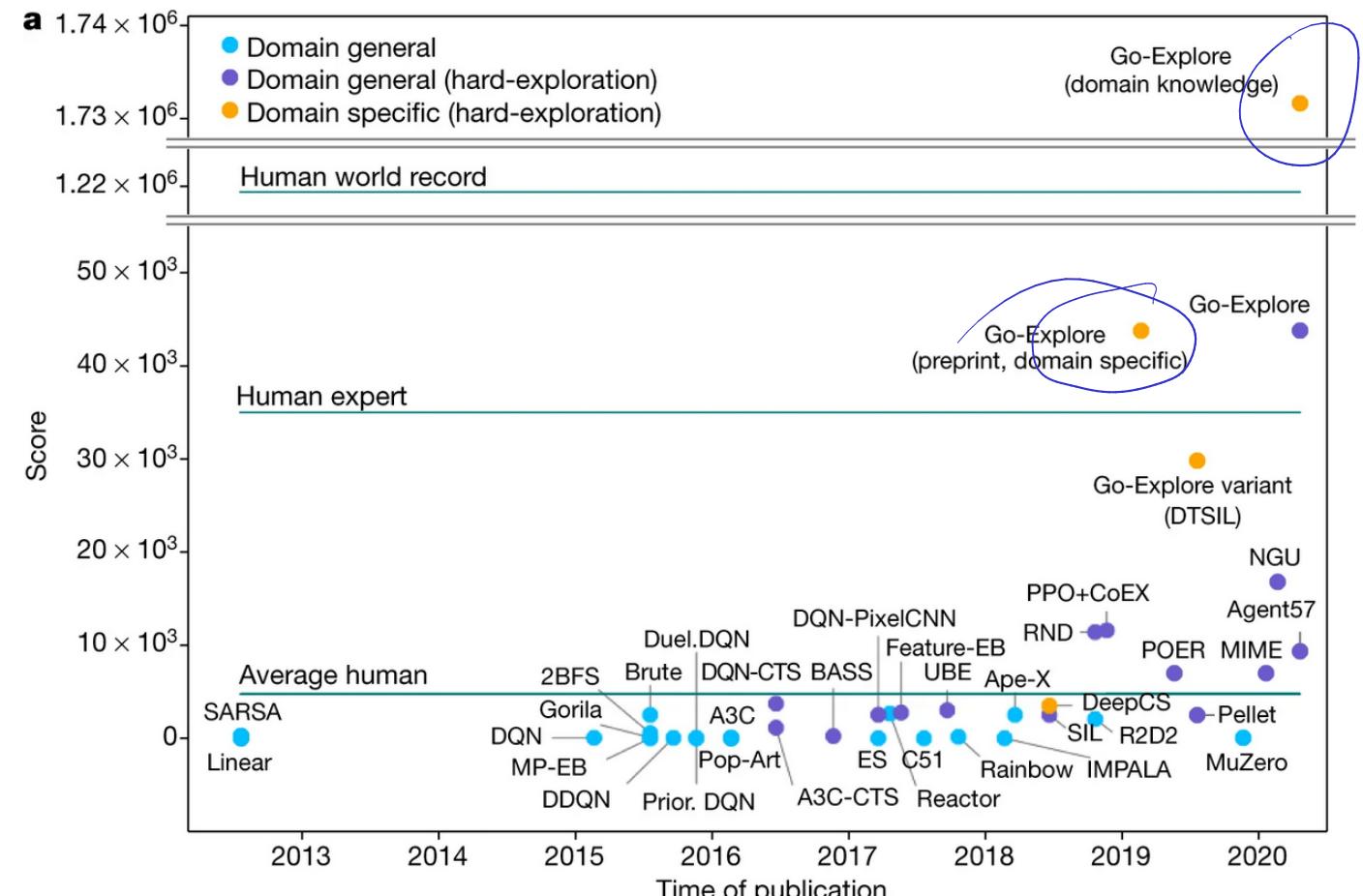
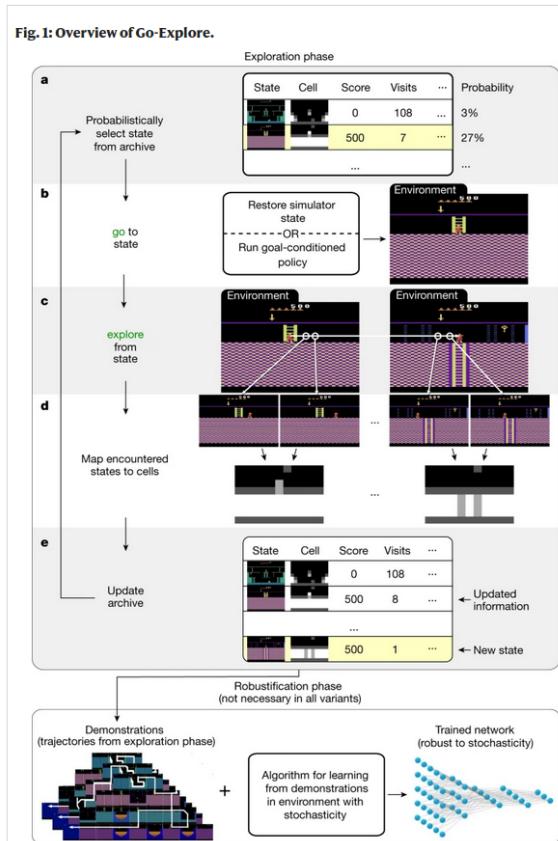


Exploration Bonus



Example 4: Go-Explore

"First return, then explore"



(Uber AI Labs)

Actor-Critic

Actor-Critic

Can we combine value-based and policy-based methods?

Actor-Critic

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

Actor-Critic

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$\overset{\uparrow}{Q(s,a)}$ $\overset{\uparrow}{V(s)}$

Actor-Critic

$$\nabla U(\theta) = E_\tau \left[\sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$\overset{\uparrow}{Q(s,a)}$ $\overset{\uparrow}{V(s)}$

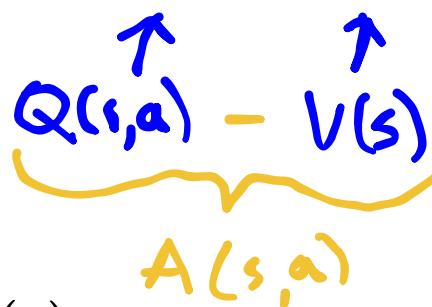
Advantage Function: $A(s, a) = Q(s, a) - V(s)$

$$\max_a Q(s, a) = V(s)$$

Actor-Critic

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$$Q(s,a) - V(s)$$



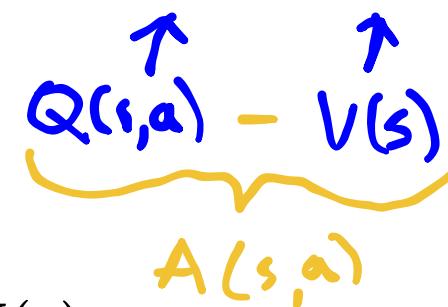
$$A(s,a)$$

Advantage Function: $A(s,a) = Q(s,a) - V(s)$

Actor-Critic

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$$Q(s,a) - V(s)$$



$$A(s,a)$$

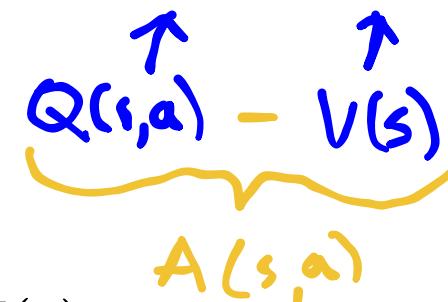
Advantage Function: $A(s, a) = Q(s, a) - V(s)$

- Actor: π_{θ}

Actor-Critic

$$\nabla U(\theta) = E_\tau \left[\sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k | s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$$Q(s,a) - V(s)$$



Advantage Function: $A(s, a) = Q(s, a) - V(s)$

- Actor: π_θ
- Critic: Q_ϕ and/or A_ϕ and/or V_ϕ

Actor-Critic

$$\nabla U(\theta) = E_\tau \left[\sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k | s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$$Q(s,a) - V(s)$$

Advantage Function: $A(s, a) = Q(s, a) - V(s)$

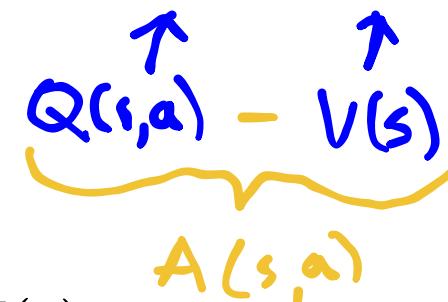
- Actor: π_θ
- Critic: Q_ϕ and/or A_ϕ and/or V_ϕ

Alternate between training Actor and Critic

Actor-Critic

$$\nabla U(\theta) = E_\tau \left[\sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k | s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]$$

$$Q(s,a) - V(s)$$



Advantage Function: $A(s, a) = Q(s, a) - V(s)$

- Actor: π_θ
- Critic: Q_ϕ and/or A_ϕ and/or V_ϕ

Alternate between training Actor and Critic

Problem: Instability

Actor-Critic

$Q(s, a)$

Which should we learn? A , Q , or V ?

$V(s)$

Actor-Critic

Which should we learn? A , Q , or V ?

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \underbrace{\gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)}_{\text{critic}}) \right]$$

Actor-Critic

Which should we learn? A , Q , or V ?

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (\underbrace{r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)}_{\text{temporal difference residual}}) \right]$$

-

Actor-Critic

Which should we learn? A , Q , or V ?

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

temporal difference residual

$$l(\phi) = E \left[(V_{\phi}(s) - V^{\pi_{\theta}}(s))^2 \right]$$

.

Actor-Critic

Which should we learn? A , Q , or V ?

$$\nabla U(\theta) = E_{\tau} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_k + \gamma V_{\phi}(s_{k+1}) - V_{\phi}(s_k)) \right]$$

temporal difference residual

$$l(\phi) = E \left[(V_{\phi}(s) - V^{\pi_{\theta}}(s))^2 \right]$$

*estimate with
reward to go
from sims*

Generalized Advantage Estimation

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k)$$

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k)$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t$$

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k)$$

V^{π_θ}

$$A(s_k, a_k) \approx \underbrace{\sum_{t=k}^{\infty} \gamma^{t-k} r_t}_{r_{\text{to go}}}$$

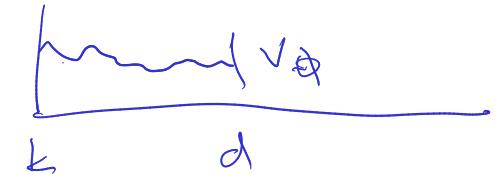
← High Bias
← High Variance

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \underbrace{\sum_{t=k}^{d-1} \gamma^{t-k} r_t}_{\text{Bias}} + \underbrace{\gamma^{d-k} r_d}_{\text{Variance}} + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$



Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

$$\text{let } \delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

$$\text{let } \delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

$$A_{\text{GAE}}(s_k, a_k) \approx \sum_{t=k}^{\infty} (\gamma \lambda)^{t-k} \delta_t$$

Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

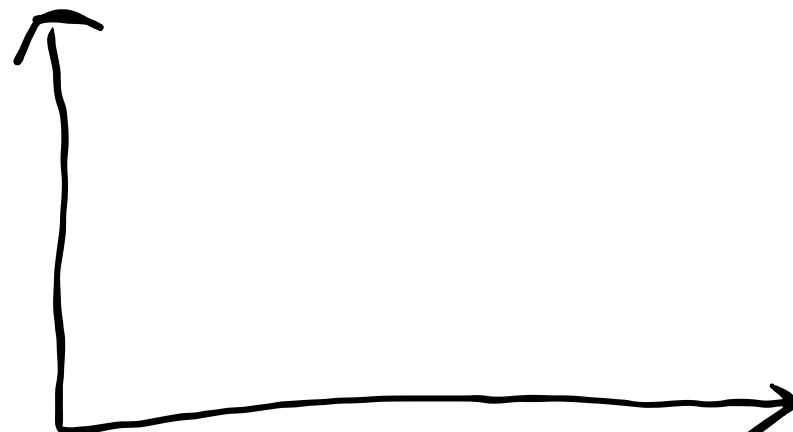
$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

let $\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$

$$A_{\text{GAE}}(s_k, a_k) \approx \sum_{t=k}^{\infty} (\gamma \lambda)^{t-k} \delta_t$$



Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k) \quad \leftarrow \text{High Bias}$$

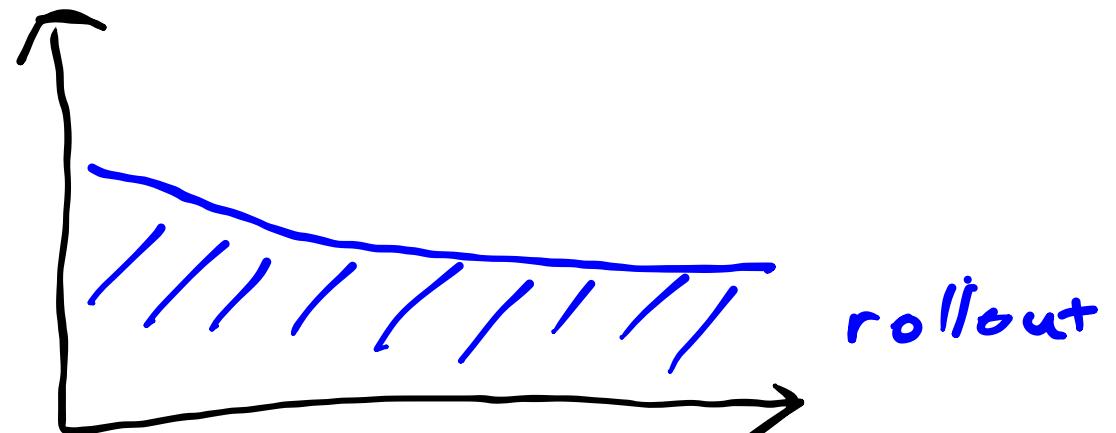
$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t \quad \leftarrow \text{High Variance}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

let $\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$

$$A_{\text{GAE}}(s_k, a_k) \approx \sum_{t=k}^{\infty} (\gamma \lambda)^{t-k} \delta_t$$



Generalized Advantage Estimation

$$A(s_k, a_k) \approx r_k + \gamma V_\phi(s_{k+1}) - V_\phi(s_k)$$

← High Bias

$$V^{\pi_\phi}$$

$$A(s_k, a_k) \approx \sum_{t=k}^{\infty} \gamma^{t-k} r_t$$

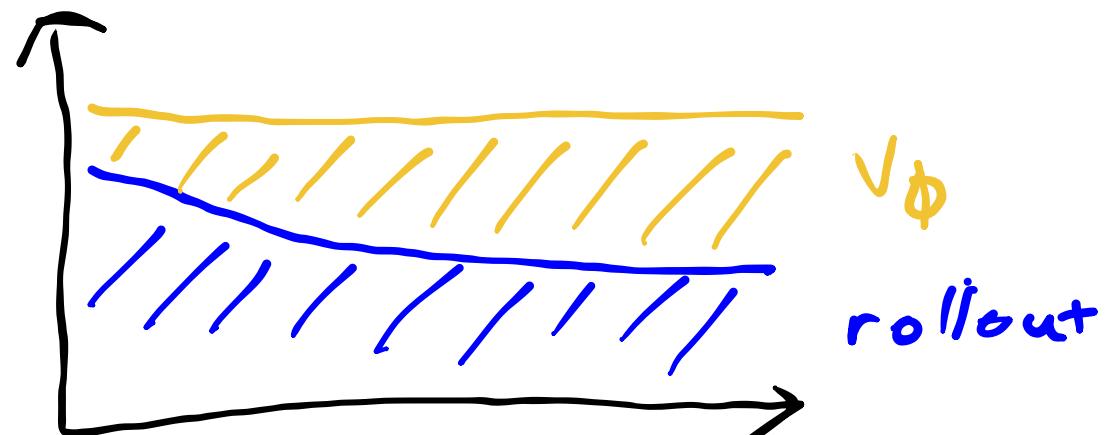
← High Variance

$$A(s_k, a_k) \approx \sum_{t=k}^{d-1} \gamma^{t-k} r_t + \gamma^{d-k} r_d + \gamma V_\phi(s_{d+1}) - V_\phi(s_d)$$

When should we stop?

let $\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$

$$A_{\text{GAE}}(s_k, a_k) \approx \sum_{t=k}^{\infty} (\gamma \lambda)^{t-k} \delta_t$$



Continuous Actions: Deep Deterministic Policy Gradient

$$Q_\phi(s, a)$$

$$\underset{\uparrow}{\operatorname{argmax}} Q_\phi$$

$$\pi_\phi(s) = \underset{a \in A}{\operatorname{argmax}} Q_\phi(s, a)$$

DDPG - fragile

$$l(\phi) = E[r + \gamma Q_\phi(s', \pi_\phi(s')) - Q_\phi(s, a)]^2]$$

$$U(\theta) = E[Q_\phi(s, \pi_\phi(s))]$$

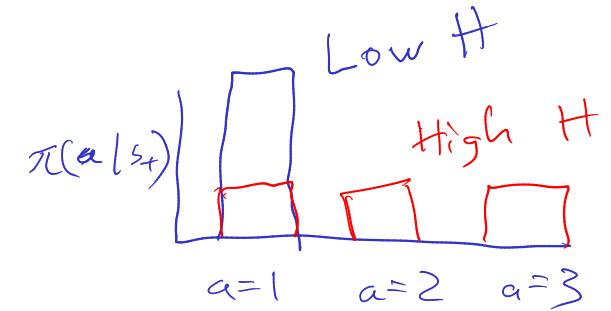
$$\begin{aligned} \nabla U(\theta) &= E[\nabla_\theta Q_\phi(s, \pi_\phi(s))] \\ &\equiv E[\nabla_\theta \pi_\phi(s) \nabla_a Q_\phi(s, a) \Big|_{a=\pi_\phi(s)}] \end{aligned}$$

Soft Actor Critic: Entropy Regularization

Soft Actor Critic: Entropy Regularization

$$\sum \gamma^t r^t$$

$$J(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \underbrace{\alpha \mathcal{H}(\pi(\cdot | s_t))}_{\text{temp}}) \right]$$



Soft Actor Critic: Entropy Regularization

$$J(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

$$V(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)]$$

Soft - Value

Soft Actor Critic: Entropy Regularization

$$J(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

$$V(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)]$$

$$\mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [\underline{V(\mathbf{s}_{t+1})}]$$

Soft Actor Critic: Entropy Regularization

$$J(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right]$$

Soft Policy Iteration

$$\begin{aligned} V(\mathbf{s}_t) &= \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)] \\ \mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t) &\triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V(\mathbf{s}_{t+1})] \\ \pi_{\text{new}} &= \arg \min_{\pi' \in \Pi} D_{\text{KL}} \left(\pi'(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot))}{Z^{\pi_{\text{old}}}(\mathbf{s}_t)} \right) \end{aligned}$$

$$\frac{e^{\lambda Q}}{\sum e^{\lambda Q}}$$

Soft Actor Critic

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

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

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

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

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

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

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

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

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

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} (Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t))^2 \right]$$

entropy Bonus

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

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

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

for each gradient step **do**

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

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

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

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

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

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

end for

end for

Soft Actor Critic

Algorithm 1 Soft Actor-Critic

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

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

for each gradient step **do**

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

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

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

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

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} (Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t))^2 \right]$$

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

end for

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

end for

Soft Actor Critic

Advantages:

Soft Actor Critic

Advantages:

- Stable

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies
- Exploration

Soft Actor Critic

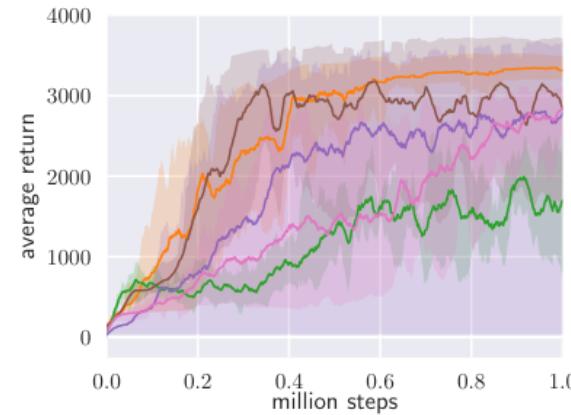
Advantages:

- Stable
- Learns many near-optimal policies
- Exploration
- Insensitivity to hyperparameters

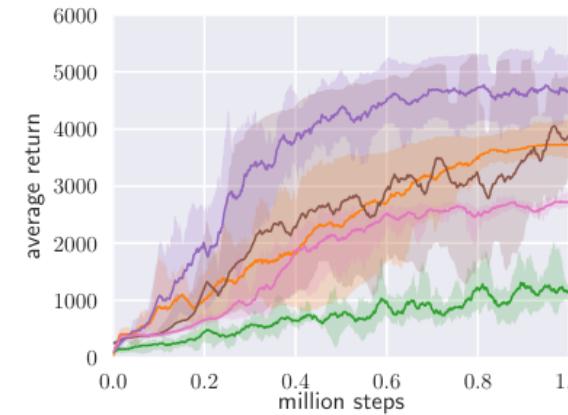
Soft Actor Critic

Advantages:

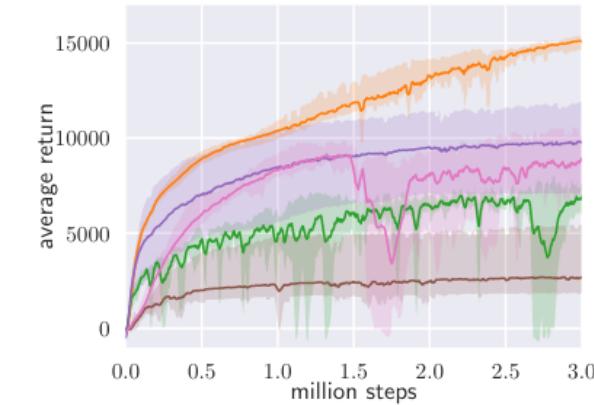
- Stable
- Learns |
- Explora
- Insensit



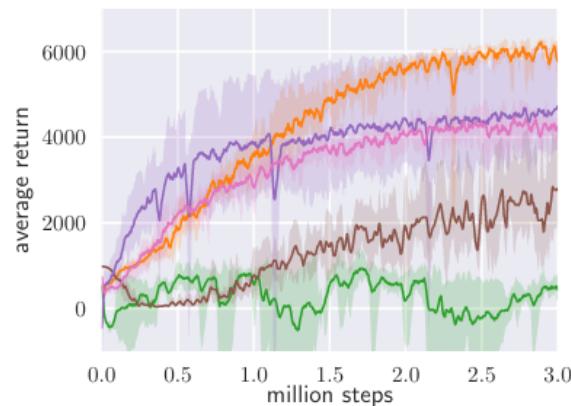
(a) Hopper-v1



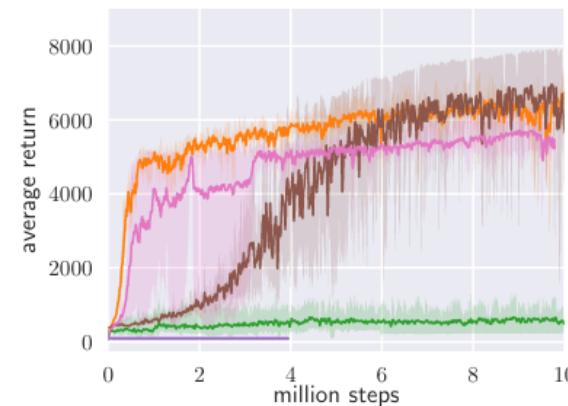
(b) Walker2d-v1



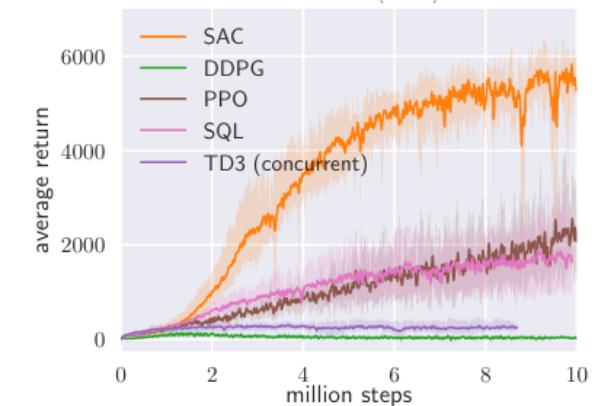
(c) HalfCheetah-v1



(d) Ant-v1



(e) Humanoid-v1



(f) Humanoid (rllab)

Soft Actor Critic

Advantages:

- Stable
- Learns many near-optimal policies
- Exploration
- Insensitivity to hyperparameters

Disadvantages

- Sensitive to α Solution = Entropy
constraint* and adjust α

$$\begin{aligned} & \text{maximize } E\left[\sum_t r_t\right] \\ \text{st} \quad & H(\pi(\cdot|s)) \geq h \end{aligned}$$

Soft Actor Critic

Advantages:

Disadvantages

- Starts with **Algorithm 1 Soft Actor-Critic**
 - Learning loop:
 - Input: θ_1, θ_2, ϕ
 - $\theta_1 \leftarrow \theta_1, \theta_2 \leftarrow \theta_2$
 - $\mathcal{D} \leftarrow \emptyset$
 - for each iteration do
 - for each environment step do
 - $a_t \sim \pi_\phi(a_t | s_t)$
 - $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$
 - $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$
 - end for
 - for each gradient step do
 - $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$
 - $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$
 - $\alpha \leftarrow \alpha - \lambda \hat{\nabla}_\alpha J(\alpha)$
 - $\theta_i \leftarrow \tau \theta_i + (1 - \tau) \theta_i$ for $i \in \{1, 2\}$
 - end for
 - end for
 - Output: θ_1, θ_2, ϕ
-
- ▷ Initial parameters
 - ▷ Initialize target network weights
 - ▷ Initialize an empty replay pool
 - ▷ Sample action from the policy
 - ▷ Sample transition from the environment
 - ▷ Store the transition in the replay pool
 - ▷ Update the Q-function parameters
 - ▷ Update policy weights
 - ▷ Adjust temperature
 - ▷ Update target network weights
 - ▷ Optimized parameters

Wisdom

- Brittleness
- Sample Complexity
- Applications (alpha fold)

