

# **FEDERATED ENSEMBLE-DIRECTED OFFLINE REINFORCEMENT LEARNING ALGORITHM (FEDORA)**

SEONVIN CHO

**INFORMATION AND INTELLIGENCE SYSTEMS LAB.**

ELECTRONIC ENGINEERING, HANYANG UNIVERSITY

**February 18, 2025**

# FEDERATED OFFLINE REINFORCEMENT LEARNING

---

- **Goal**

- To learn the optimal policy using only offline data from the operational policies of **multiple clients** with **different levels of expertise**
- **Without** the clients knowing the **quality** of their data, or sharing it with one another or the server

- **Challenges**

- Ensemble Heterogeneity: Learn policies of varying quality
- Pessimistic Value Computation: Q-value underestimation due to limited client datasets
- Data Heterogeneity: Varying data quality

⇒ **Federated Ensemble-Directed Offline RL Algorithm(FEDORA)**

# RELATED WORK

---

- **Federated Learning**

- To minimize  $F(\theta) = \mathbb{E}_{i \sim P}[F_i(\theta)]$ .
- **FedAvg algorithm:**  $\theta^{t+1} = \sum_{i=1}^{|N|} \omega_i \theta_i^t$ , where  $\omega_i = \frac{|D_i|}{\sum_{j=1}^{|N|} |D_j|}$ .

- **Reinforcement Learning**

- To maximize  $J(\pi) = \mathbb{E}_{\pi, P, \mu}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ .

- **Offline Reinforcement Learning**

- To learn  $\pi$  only using a static dataset by  $\pi_b$  without any additional interactions with the environment.
- Utilize the regularization to prevent distribution shift.

# RELATED WORK

- **Federated Learning**

- To minimize  $F(\theta) = \mathbb{E}_{i \sim P}[F_i(\theta)]$ .

- **FedAvg algorithm:**  $\theta^{t+1} = \sum_{i=1}^{|N|} \omega_i \theta_i^t$ , where  $\omega_i = \frac{|D_i|}{\sum_{j=1}^{|N|} |D_j|}$ .

- **Reinforcement Learning**

- To maximize  $J(\pi) = \mathbb{E}_{\pi, P, \mu}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ .

- **Offline Reinforcement**

- To learn  $\pi$  only using a static dataset.

- Each client learns using its own dataset under specific behavior policies.
    - The learned policy varies depending on the behavior policy.
    - Simply aggregate all client models degrades performance.

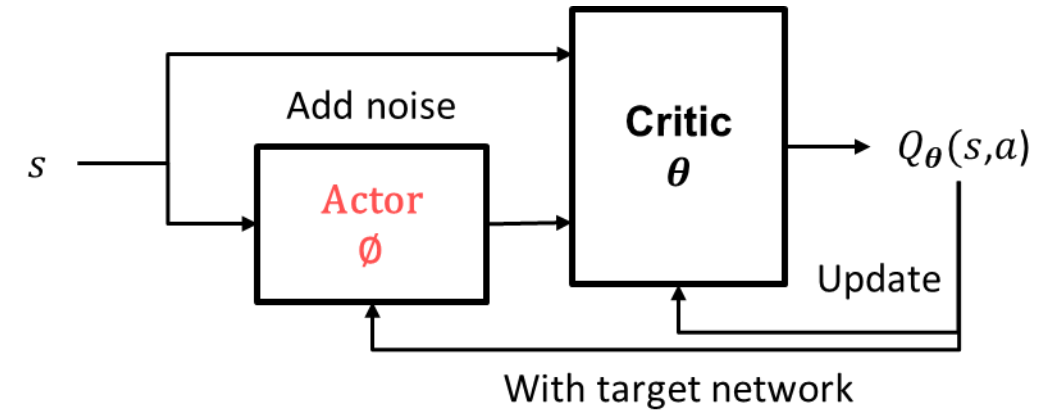
- Utilize the regularization to prevent distribution shift.

- **TD3-BC**(Twin Delayed DDPG-Behavior Cloning): To prevent distribution shift.

# TD3-BC

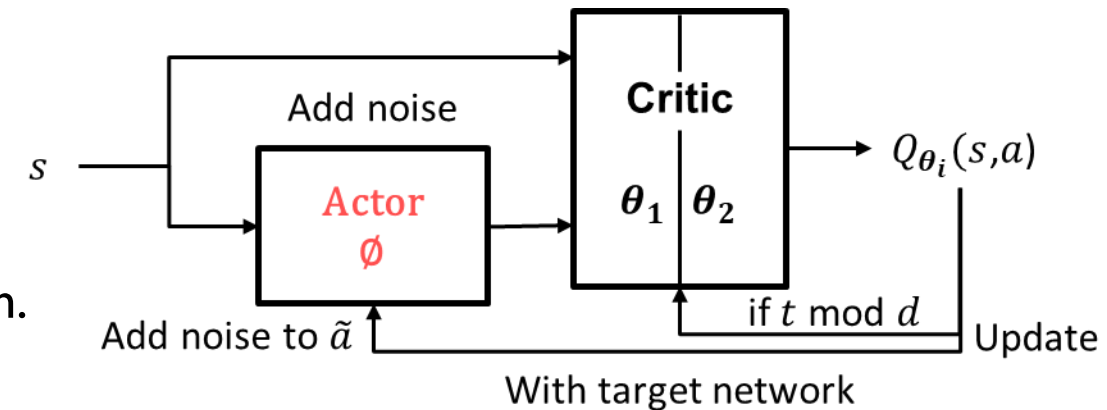
- DDPG

- Overestimate Q-values in critic.
- Update actor at every step  $\rightarrow$  instability in Q-value.
- Change target Q-values too rapidly  $\rightarrow$  unstable.



- TD3(Twin Delayed DDPG)

- Utilize twin Q-networks( $Q_{\theta_1}, Q_{\theta_2}$ ) and update minimum.
- Delay updating actor compared to critic.
- Smooth target policy with adding gaussian noise to action when computing target Q-value.



# TD3-BC

- DDPG

- Overestimate Q-values in critic.
- Update actor at every step  $\rightarrow$  instability in Q-value.
- Change target Q-values too rapidly  $\rightarrow$  unstable.

- TD3(Twin Delayed DDPG)

- Utilize twin Q-networks( $Q_{\theta_1}, Q_{\theta_2}$ ) and update minimum.
- Delay updating actor compared to critic.
- Smooth target policy with adding gaussian noise to action when computing target Q-value.

## DDPG

- $y \leftarrow r + \gamma Q_{\theta'}(s', \tilde{a})$
- Update  $\phi$  (w.r.t. actor policy  $\pi_{\phi}$ ) at every  $t$
- $\tilde{a} \leftarrow \pi_{\phi'}(s')$

## TD3

- $y \leftarrow r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tilde{a})$
- Update  $\phi$  (w.r.t. actor policy  $\pi_{\phi}$ ) if  $t \bmod d$
- $\tilde{a} \leftarrow \pi_{\phi'}(s') + \epsilon$  where  $\epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$

# TD3-BC

---

- TD3 Problems

- Not suitable for offline RL. ( $\because$  requires exploration through interactions with the environment)
- The actor may select actions that deviate from the original data distribution in pursuit of optimal Q.

- TD3-BC

- Add BC regulation term
  - › To favor actions contained in the dataset  $\mathcal{D}$
  - › To use only original data without exploration

$$\pi \leftarrow \operatorname{argmax}_{\pi} \mathbb{E}_{s \sim \mathcal{D}} [\lambda Q(s, \pi(s)) - (\pi(s) - a)^2]$$

- TD3 term:  $Q(s, \pi(s))$  for maximize Q-value.
- BC term:  $-(\pi(s) - a)^2$  for reducing the difference between action and policy.

# FEDORA

- **Solution**

- Ensemble Heterogeneity → Ensemble-directed learning to weigh client contribution.

- › Weights ~ entropy regularization  $\omega_i = \frac{e^{\beta J_i |D_i|}}{\sum_j e^{\beta J_j |D_j|}}$  where  $J_i^t = \mathbb{E}_{s \sim D_i}[Q_i^t(s, \pi_i^t(s))]$ .

- › Federated policy ~ weighted combination of client policies  $\pi_{fed}^{t+1} = \sum_i \omega_i \pi_i^t$ .

- Pessimistic Value Computation → Federated optimism for critic training.

- › Ensemble-directed Federation → Optimistic target  $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ .

- Data Heterogeneity → Proximal policy update.

- ›  $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{actor}(\pi)$  where  $\mathcal{L}_{actor}(\pi) = \mathcal{L}_{local}(\pi) + \mathbb{E}_{(s,a) \sim D_i}[(\pi(s) - \pi_{fed}^{t+1})^2]$ ,

- ›  $\mathcal{L}_{local}(\pi) = \mathbb{E}_{(s,a) \sim D_i}[-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$



# FEDORA

- **Solution**

- Ensemble Heterogeneity → Ensemble-directed learning to weigh client contribution.

- › Weights ~ entropy regularization  $\omega_i = \frac{e^{\beta J_i |D_i|}}{\sum_j e^{\beta J_j |D_j|}}$  where  $J_i^t = \mathbb{E}_{s \sim D_i}[Q_i^t(s, \pi_i^t(s))]$ .

- › Federated policy ~ weighted combination of client policies  $\pi_{fed}^{t+1} = \sum_i \omega_i \pi_i^t$ .

- Pessimistic Value Computation → Federated optimism for critic training.

- › Ensemble-directed Federation → Optimistic target  $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ .

- Data Heterogeneity → Proximal policy update.

- ›  $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{actor}(\pi)$ 

TD3-BC:  $\pi \leftarrow \operatorname{argmax}_{\pi} \mathbb{E}_{s \sim \mathcal{D}}[\lambda Q(s, \pi(s)) - (\pi(s) - a)^2]$

FEDORA:  $\pi \leftarrow \operatorname{argmin}_{\pi} \mathbb{E}_{(s,a) \sim \mathcal{D}}[-\lambda Q(s, \pi(s)) + (\pi(s) - a)^2]$

- ›  $\mathcal{L}_{local}(\pi) = \mathbb{E}_{(s,a) \sim D_i}[-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$

# FEDORA

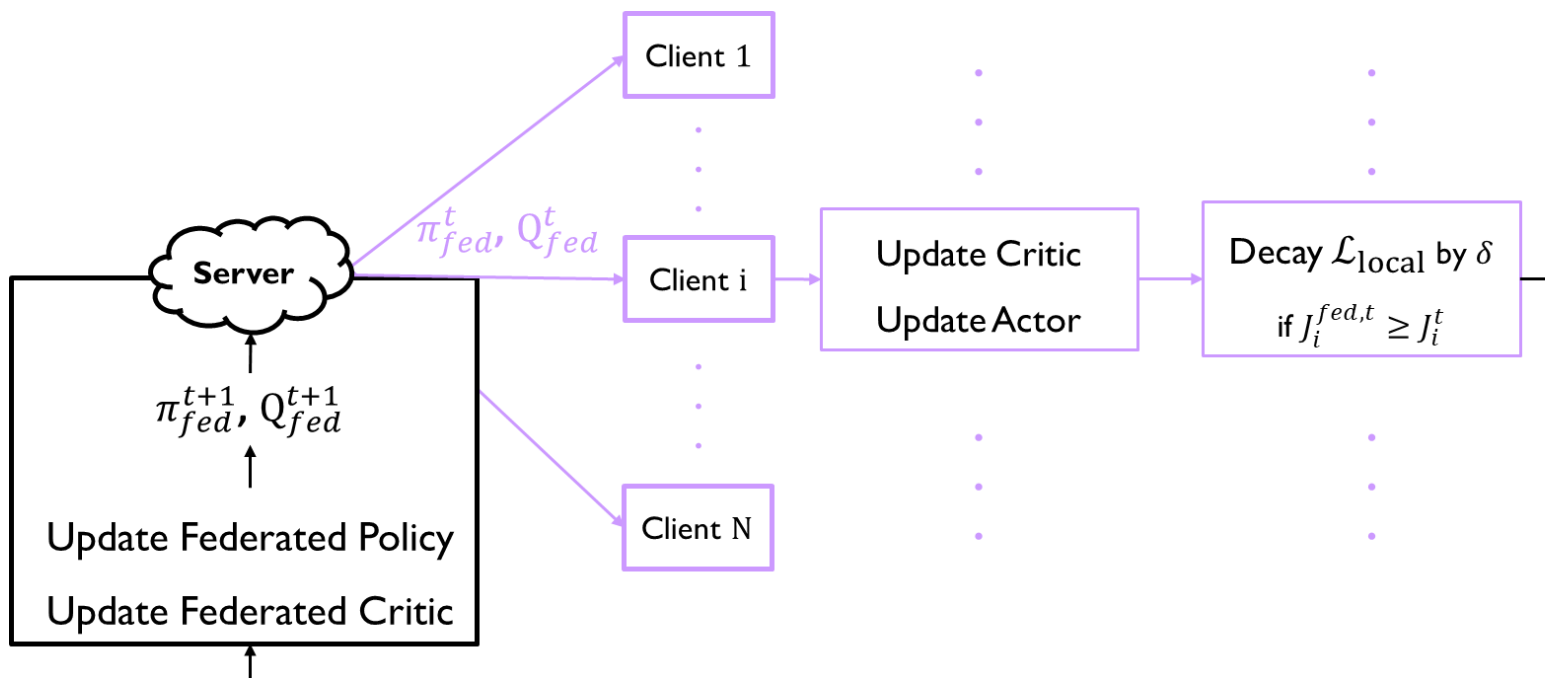
$$\text{Eq. (8)} \quad \omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}, \quad \pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$$

$$\text{Eq. (9)} \quad Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$$

$$\text{Eq. (10)} \quad Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$$

$$\text{Eq. (11)} \quad \mathcal{L}_{actor}(\pi) = \mathcal{L}_{local}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2],$$

$$\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{actor}(\pi)$$



## Algorithm 1 Outline of Client $i$ 's Algorithm

- 1: **function** train\_client( $\pi_{fed}^t, Q_{fed}^t$ )
- 2:  $\pi_i^{(t,0)} = \pi_{fed}^t, \quad Q_i^{(t,0)} = Q_{fed}^t$
- 3: **for**  $1 \leq k < K$  **do**
- 4:     Update Critic by one gradient step w.r.t. Eq. (10)
- 5:     Update Actor by one gradient step w.r.t. Eq. (11)
- 6: **end for**
- 7:     Decay  $\mathcal{L}_{local}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$
- 8: **end function**

## Algorithm 2 Outline of Server Algorithm

- 1: Initialize  $\pi_{fed}^1, Q_{fed}^1$
- 2: **for**  $t \in 1 \dots$  **do**
- 3:     Send  $\pi_{fed}^t$  and  $Q_{fed}^t$  to  $i \in \mathcal{N}$
- 4:     Sample  $\mathcal{N}_t \subset \mathcal{N}$
- 5:     **for**  $i \in \mathcal{N}_t$  **do**
- 6:          $i.train\_client(\pi_{fed}^t, Q_{fed}^t)$  (Client side)
- 7:     **end for**
- 8:     Compute  $\pi_{fed}^{t+1}$  and  $Q_{fed}^{t+1}$  for clients in  $\mathcal{N}_t$  using Eq. (8) and (9) respectively.
- 9: **end for**

# FEDORA OVERVIEW

Server

```

Initialize  $\pi_{fed}^1, Q_{fed}^1$ 
for  $t \in 1, \dots$  (# of round) do
    Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$ 
    Sample  $N_t \subset N$ 
    for  $i \in N_t$  do
        train 'i' th client
    endfor
    Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$ 
    Ensemble  $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$ 
    Federation  $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$ 
endfor
    
```

$\pi_{fed}^t$   
 $Q_{fed}^t$

$\omega_i^t$   
 $\pi_i^t$   
 $Q_i^t$

Client

```

 $\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$ 
for  $1 \leq k < K$  do
    Federated  $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ 
    Optimism  $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$ 

    Proximal  $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$ 
    Policy  $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$ 
     $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$ 
endfor

 $J_i^{fed,t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$ 
Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$ 

 $\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$ 
    
```

# FEDORA OVERVIEW

Server

```

Initialize  $\pi_{fed}^1, Q_{fed}^1$ 
for  $t \in 1, \dots$  (# of round) do
    Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$ 
    Sample  $N_t \subset N$ 
    for  $i \in N_t$  do
        train 'i' th client
    endfor
    Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$ 
     $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$ 
     $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$ 
endfor
    
```

$\pi_{fed}^t$   
 $Q_{fed}^t$

$\omega_i^t$   
 $\pi_i^t$   
 $Q_i^t$

Pessimistic Value Computation → Federated optimism for critic training.

Optimistic target  $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ .

Client

```

 $\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$ 
for  $1 \leq k \leq K$  do
    
```

Update Critic

$\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$  Federated Optimism  
 $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$

$\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$

$\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$

$\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$

endfor

$J_i^{\text{fed},t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$

Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{\text{fed},t} \geq J_i^t$

$\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$

# FEDORA OVERVIEW

Server

```

Initialize  $\pi_{fed}^1, Q_{fed}^1$ 
for  $t \in 1, \dots$  (# of round) do
    Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$ 
    Sample  $N_t \subset N$ 
    for  $i \in N_t$  do
        train 'i' th client
    endfor
    Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$ 
     $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$ 
     $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$ 
endfor
    
```

$\pi_{fed}^t$   
 $Q_{fed}^t$

$\omega_i^t$   
 $\pi_i^t$   
 $Q_i^t$

Client

```

 $\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$ 
for  $1 \leq k < K$  do
     $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ 
     $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$ 

Regularize local offline data
         $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$ 
        Regularize federated policy
         $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$ 
         $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$


endfor
    
```

$J_i^{fed,t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$  **Update Actor**  
 Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$  **Proximal Policy Update**

$$\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$$

Data Heterogeneity → Proximal policy update.

›  $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$  where  $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$ ,

›  $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$

# FEDORA OVERVIEW

Server

```

Initialize  $\pi_{fed}^1, Q_{fed}^1$ 
for  $t \in 1, \dots$  (# of round) do
    Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$ 
    Sample  $N_t \subset N$ 
    for  $i \in N_t$  do
        train 'i' th client
    endfor
    Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$ 
     $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$ 
     $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$ 
endfor
    
```

$\pi_{fed}^t$   
 $Q_{fed}^t$

$\omega_i^t$   
 $\pi_i^t$   
 $Q_i^t$

Client

```

 $\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$ 
for  $1 \leq k < K$  do
     $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ 
     $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$ 

     $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$ 
     $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$ 
     $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$ 
endfor

 $J_i^{fed,t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$ 
Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$ 

 $\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$ 
    
```

Local estimate

Updated critic

Decay the influence of local data

# FEDORA OVERVIEW

Server

Initialize  $\pi_{fed}^1, Q_{fed}^1$   
 for  $t \in 1, \dots$  (# of round) do  
     Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$   
     Sample  $N_t \subset N$   
     for  $i \in N_t$  do  
         train 'i' th client  
     endfor  
     Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$   
     Update Policy  $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$   
     Update Critic  $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$   
 endfor

## Ensemble-Directed Federated Learning

Ensemble Heterogeneity → Ensemble-directed learning to weigh client contribution.

- Weights ~ entropy regularization  $\omega_i = \frac{e^{\beta J_i^t |D_i|}}{\sum_j e^{\beta J_j^t |D_j|}}$  where  $J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$ .
- Federated policy ~ weighted combination of client policies  $\pi_{fed}^{t+1} = \sum_i \omega_i \pi_i^t$ .

Pessimistic Value Computation → Federated optimism for critic training.

- Optimistic target  $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ .

Client

$\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$   
 for  $1 \leq k < K$  do  
      $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$   
      $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$   
  
      $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$   
      $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$   
      $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$   
 endfor

Local estimate

$$J_i^{fed,t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$$

Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$

Updated critic

$$\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$$

# FEDORA OVERVIEW

Server

```

Initialize  $\pi_{fed}^1, Q_{fed}^1$ 
for  $t \in 1, \dots$  (# of round) do
    Send  $\pi_{fed}^t, Q_{fed}^t$  to client  $i \in N$ 
    Sample  $N_t \subset N$ 
    for  $i \in N_t$  do
        train 'i' th client
    endfor
    Compute  $\pi_{fed}^{t+1}, Q_{fed}^{t+1}$  for clients in  $N_t$ 
     $\pi_{fed}^{t+1} = \sum_{i=1}^{|N|} \omega_i^t \pi_i^t$ 
     $Q_{fed}^{t+1} = \sum_i \omega_i^t Q_i^t$ 
endfor
    
```

$\pi_{fed}^t$   
 $Q_{fed}^t$

$\omega_i^t$   
 $\pi_i^t$   
 $Q_i^t$

Client

```

 $\pi_i^{(t,0)} = \pi_{fed}^t, Q_i^{(t,0)} = Q_{fed}^t$ 
for  $1 \leq k < K$  do
     $\tilde{Q}_i^{(t,k)}(s, a) = \max(Q_i^{(t,k)}(s, a), Q_{fed}^t(s, a))$ 
     $Q_i^{(t,k+1)} = \operatorname{argmin}_Q \mathbb{E}_{(s,a,r,s') \sim D_i} [(r + \gamma \tilde{Q}_i^{(t,k)}(s', a') - Q(s, a))^2]$ 

     $\mathcal{L}_{\text{local}}(\pi) = \mathbb{E}_{(s,a) \sim D_i} [-Q_i^{(t,k)}(s, \pi(s)) + (\pi(s) - a)^2]$ 
     $\mathcal{L}_{\text{actor}}(\pi) = \mathcal{L}_{\text{local}}(\pi) + \mathbb{E}_{(s,a) \sim D_i} [(\pi(s) - \pi_{fed}^{t+1})^2]$ 
     $\pi_i^{t,k+1} = \operatorname{argmin}_{\pi} \mathcal{L}_{\text{actor}}(\pi)$ 
endfor

 $J_i^{fed,t} = \mathbb{E}_{s \sim D_i} [Q_{fed}^t(s, \pi_{fed}^t)], J_i^t = \mathbb{E}_{s \sim D_i} [Q_i^t(s, \pi_i^t(s))]$ 
Decay  $\mathcal{L}_{\text{local}}$  by  $\delta$  if  $J_i^{fed,t} \geq J_i^t$ 

 $\omega_i^t = \frac{e^{\beta J_i^t |D_i|}}{\sum_{j=1}^{|N|} e^{\beta J_j^t |D_j|}}$ 
    
```