
Algorithm 1 A Stealth and Defensive Backdoor based on Steganographic Algorithm in Federated Learning

Input: start epoch E_s , attack num E_a , end epoch E_e , client set C , selected client set C_m , adversary set C_{adv} , global model G , local model θ , central server C_s , aggregate algorithm *PartFedAvg*, benign datasets \hat{D} , poisoned datasets \hat{D}_p , benign learning rate η_b , poison learning rate η_p , PartFedAvg gradient removal scale $\mathcal{R}\%$

Output: a global model with high accuracy, stealth and defensive backdoor and high accuracy in main-task

- 1: C_s select n clients by random into C_m
 - 2: C_s build a global model G
 - 3: C_s send G to each client in C_m
 - 4: **for** epoch $< E_e$ and epoch $> E_s + E_a$ **do**
 - 5: **for** number k of client in C_m **do**
 - 6: **if** client $e_i \in C_{adv}$ **then**
 - 7: Download G as local model L and train L by private benign
 - 8: Compute gradient by \hat{D}_p on batch B_i of size ℓ
 - 9: $g_i^p = \frac{1}{\ell} \sum_{i=1}^n \nabla_{\theta} \mathcal{L}(\theta_{e_i}, \hat{D}_p)$
 - 10: Update $\theta_{e_{i+1}} = \theta_{e_i} - \eta_p g_i^p$ where $g_i^p \not\in \text{top}_{5\%}(g)$
 - 11: Upload $\theta_{e_{i+1}}$ to C_s
 - 12: **else if** client $e_i \notin C_{adv}$ **then**
 - 13: Download G as local model L and train L by private poisoned dataset
 - 14: Compute gradient by \hat{D}_b on batch B_i of size ℓ
 - 15: $g_i^b = \frac{1}{\ell} \sum_{i=1}^n \nabla_{\theta} \mathcal{L}(\theta_{e_i}, \hat{D})$
 - 16: Update $\theta_{e_{i+1}} = \theta_{e_i} - \eta_b g_i^b$
 - 17: Upload $\theta_{e_{i+1}}$ to C_s
 - 18: C_s receive $\sum_1^k \theta_{e_{i+1}}^k$ and generate update gradient U for G
 - 19: set $\mathcal{R}\%$ of gradient U to zero
 - 20: $G_{i+1} = G_i - U_i$
 - 21: **return** Final global model G with backdoor
-