
Algorithm 1 Multi-stage Reconstruction-based Membership Inference Attack

1: **Input:** Private dataset D_i for each client, initialized model ω , teacher model T , number of clients N , global rounds T
2: **Output:** Robust global model
3: **for** $t = 1, 2, \dots, T$ **do**
4: **for** $i = 1, \dots, N$ in parallel **do**
5: Send global model ω' to local client i
6: $\omega' \leftarrow \text{LocalUpdate}(\omega')$
7: **end for**
8: $L(\omega) \leftarrow \sum_{i \in [N]} \frac{D_i}{\sum_{i \in [N]} D_i} L_i(\omega') \quad (3)$
9: **end for**
10: **LocalUpdate**(ω')
11: **for** each local epoch **do**
12: **for** each batch (x_i, y_i) of D_i **do**
13: /* Adversarial examples generation */
14: $x_i^{adv} \leftarrow x_i + \delta \quad (6)$
15: /* Clean examples augmentation */
16: $x_{ij} \leftarrow \lambda x_i + (1 - \lambda)x_j \quad (8)$
17: /* Adversarial examples augmentation */
18: $x_i^{adv} \leftarrow x_{ij} + \lambda(1 - \lambda)x_i^{adv} \quad (12)$
19: /* Vanilla mixture knowledge distillation */
20: $L_{VKD} \leftarrow KL(z_{ij}, z_{ij}^{adv}) + KL(z_{ij}, z_{isj}) \quad (11)$
21: /* Adversarial mixture knowledge distillation */
22: $L_{AKD} \leftarrow KL(z_i^{adv}, z_s^{adv}) + KL(z_{ij}^{adv}, z_s^{adv}) \quad (14)$
23: /* Consistency regularization */
24: $L_{ALG} \leftarrow \lambda_{adv} \|z_s^{adv} - z_g^{adv}\|^2 \quad (15)$
25: /* Overall local objective for each client */
26: $L \leftarrow \alpha L_{VKD} + (1 - \alpha)L_{AKD} + \lambda L_{ALG} \quad (16)$
27: **end for**
28: **end for**
29: **return** ω_i
