

# Aggregation Rules of Defenses in Federated Learning



Xiaobing Chen  
Jul 7, 2021



# Outlines

## 1. Introduction

Attack schemes;  
Capabilities of adversary;  
Defenses based on aggregation rules;

## 2. Defense Methods based on Altering Aggregation Rules

Method 1: ARFL[1];  
Method 2: FL-IOWA-DQ[2];

## 3. Extra Experiments

Compare ARFL, FL-IOWA-DQ in the same setting;  
Test their performance in extreme non-i.i.d distribution;

---

1

# Introduction

---

*Federated learning (FL) is a machine learning setting where many clients collaboratively train a model under the orchestration of a central server, while keeping the training data decentralized.*

“



# Introduction

## Attack Risks of FL

- Vulnerable spots: the distributed nature, architectural design, and data constraints ;
- Adversarial Attacks: attacks on **model performance** or on data inference

## Attack Schemes

Training-time attacks(poisoning)	<b>Data poisoning*</b> <ul style="list-style-type: none"><li>• the adversary alters the client datasets used to train the model</li></ul>
	Model update poisoning <ul style="list-style-type: none"><li>• the adversary alters model updates sent to the server.</li></ul>
Inference-time attacks(evasion)	model evasion attacks <ul style="list-style-type: none"><li>• the adversary alters the data used at inference-time</li></ul>

**Data poisoning\*** will be mainly discussed in the further context.



# Introduction

## Capabilities of adversaries

<b>Model inspection</b> <ul style="list-style-type: none"><li>Whether the adversary can observe the model parameters.</li></ul>	<ul style="list-style-type: none"><li>Black box: the adversary has no ability to inspect the parameters of the mode.</li><li>Stale white box: the adversary can only observe the model while participating in the global aggregation.</li><li>White box: the adversary can observe the model all the time.</li></ul>
<b>Participant collusion</b> <ul style="list-style-type: none"><li>Whether multiple adversaries can coordinate an attack.</li></ul>	<ul style="list-style-type: none"><li>Non-colluding: there is no capability for participants to coordinate an attack.</li><li>Cross-update collusion: past client participants can coordinate with future participants on attacks to future updates to the global model.</li><li>Within-update collusion: current client participants can coordinate on an attack to the current model update.</li></ul>
<b>Adaptability</b> <ul style="list-style-type: none"><li>Whether an adversary can alter the attack parameters as the attack progresses.</li></ul>	<ul style="list-style-type: none"><li>Static: the adversary must fix the attack parameters at the start of the attack and cannot change them.</li><li>Dynamic: the adversary can adapt the attack as training progresses.</li></ul>

Stale white box, non-colluding and static are considered in further context.



# Introduction

---

## Defense schemes

- Many existing defense methods in distributed datacenter learning and centralized learning are hard to be deployed in federated learning, such as, data sanitization and network pruning.
- Designing robust aggregators by replacing the averaging aggregation in the server is one direction that has been widely explored.
- Dynamic Quantifier federated aggregation operator (FL-IOWA-DQ) and Auto-weighted Robust Federated Learning (ARFL) will be introduced.

2

## **Defense Methods based on Altering Aggregation Rules**





## Aggregation of FedAvg

TABLE I: Summary of main notations

$N$	total number of the clients
$\bar{L}_i$	the local training loss of ith client
$D_i$	ith local training dataset
$m_i$	the size of ith local training dataset
$M$	the size of total training dataset
$\omega_t$	the global weight in time t
$\alpha_i$	weight assigned to ith local update
$\delta_i(t)$	local update of ith client in time t

FedAvg:

$$\begin{aligned}\omega_{t+1} &= \omega_t + \sum_{i=1}^N \alpha_i \delta_i(t) \\ \delta_i(t) &= \omega_t(D_i) - \omega_t\end{aligned}$$

$$\begin{aligned}\alpha_i &= \frac{m_i}{M} \\ s.t. \quad \alpha &\in \mathbb{R}_+^n, \mathbf{1}^\top \alpha = 1,\end{aligned}$$

**Larger local datasets, higher weights.**

2.1

## *Auto-weighted Robust Federated Learning (ARFL)*

---



## Aggregation of **ARFL**

$$\min_{\mathbf{w}, \alpha} \sum_{i=1}^N \alpha_i \hat{\mathcal{L}}_i(\mathbf{w}) \quad \longrightarrow \quad \min_{\mathbf{w}, \alpha} \sum_{i=1}^N \alpha_i \hat{\mathcal{L}}_i(\mathbf{w}) + \frac{\lambda}{2} \sum_{i=1}^N \frac{\alpha_i^2}{m_i}$$

Objective for FedAvg

Objective for **ARFL**

- Add **L2 regularization** based on data size



## Aggregation of ARFL

**Theorem 2.** For any  $\mathbf{w}$ , when  $\lambda > 0$  and  $\{\hat{\mathcal{L}}_i(\mathbf{w})\}_{i=1}^N$  are sorted in increasing order:  $\hat{\mathcal{L}}_1(\mathbf{w}) \leq \hat{\mathcal{L}}_2(\mathbf{w}) \leq \dots \leq \hat{\mathcal{L}}_N(\mathbf{w})$ , by setting:

$$p = \operatorname{argmax}_k \left\{ 1 + \frac{M_k(\bar{\mathcal{L}}_k(\mathbf{w}) - \hat{\mathcal{L}}_k(\mathbf{w}))}{\lambda} > 0 \right\}, \quad (7)$$

where  $M_k = \sum_{i=1}^k m_i$ ,

$$\bar{\mathcal{L}}_k(h) = \frac{\sum_{i=1}^k m_i \hat{\mathcal{L}}_i(\mathbf{w})}{M_k} \quad (8)$$

is the average loss over the first  $k$  clients that have the smallest empirical risks. Then the optimal  $\alpha$  to the problem (6) is given by:

$$\alpha_i(\mathbf{w}) = \frac{m_i}{M_p} \left[ 1 + \frac{M_p(\bar{\mathcal{L}}_p(\mathbf{w}) - \hat{\mathcal{L}}_i(\mathbf{w}))}{\lambda} \right]_+, \quad (9)$$

where  $[\cdot]_+ = \max(0, \cdot)$ .

**Lower training loss, higher weights.**



# Aggregation of ARFL

---

**Algorithm 1** Optimization of ARFL

---

**Server executes:**

```
1: Initialize  $w_0, \hat{\mathcal{L}}, \alpha$ 
2: for each round  $t = 1, 2, \dots$  do
3:   Select a subset  $S_t$  from  $N$  clients at random
4:   Broadcast the global model  $w_t$  to selected clients  $S_t$ 
5:   for each client  $i \in S_t$  in parallel do
6:      $w_{t+1}^i, \hat{\mathcal{L}}_i \leftarrow \text{ClientUpdate}(i, w_t)$ 
7:   end for
8:   Update  $w_{t+1}$  according to Eq. (11)
9:   Update  $\alpha$  according to Theorem 2
10: end for
11:
```

**ClientUpdate( $i, w$ ):** // Run on client  $i$

```
12:  $\mathcal{L}_i \leftarrow$  (evaluate training loss using training set)
13:  $\mathcal{B} \leftarrow$  (split local training set into batches of size  $B$ )
14: for each local epoch  $i$  from 1 to  $E$  do
15:   for batch  $b \in \mathcal{B}$  do
16:      $w \leftarrow w - \eta \nabla \ell(w; b)$ 
17:   end for
18: end for
19: return  $w$  and  $\hat{\mathcal{L}}_i$ 
```

---

- Line 9: weights are maintained by the server; only updates the losses from those selected clients while keeping the others unchanged
- Line 12: clients submit the training losses, together with updates



## Experiments: Datasets

	CIFAR-10[1]	FEMNIST[2]	Shakespeare
#Classes:	10;	62;	80;
#Clients:	100;	1039;	71;
#Samples:	60,000;	236,500;	417,469;
i.i.d:	Yes;	No;	No;
Model:	CNN;	CNN;	LSTM;
Task	Image	Image	Next-character
#Participants(ratio)	classification 20(20%);	classification 32(3%);	prediction 16(22%);

**i.i.d: “each local client has approximately the same amount of samples and in proportion to each of the classes. ”**

**Non i.i.d in source: each speaking role or writer is treated as a client.**

[1]. Learning Multiple Layers of Features from Tiny Images, Alex Krizhevsky, 2009.

[2]. Caldas, Sebastian, et al. "Leaf: A benchmark for federated settings." arXiv preprint arXiv:1812.01097 (2018).



## Experiments: Models

- **FedAvg**: The standard Federated Averaging aggregation approach that just calculates the weighted average of the parameters from local clients.
- **RFA[1]**: A robust aggregation approach that minimizes the weighted Geometric Median(GM) of the parameters from local clients.
- **MKrum (Multi-Krum)[2]**: A Byzantine tolerant aggregation rule, which computes a distance-related score for each update, and then add the averaged updates with high scores to the model.
- **CFL[3]**: A Clustered Federated Learning (CFL) approach that separates the client population into different groups based on the pairwise cosine similarities between their parameter updates. Updates in benign group are aggregated.

[1]. Pillutla, Krishna, Sham M. Kakade, and Zaid Harchaoui. "Robust aggregation for federated learning." arXiv preprint arXiv:1912.13445 (2019).

[2]. Blanchard, Peva, et al. "Machine learning with adversaries: Byzantine tolerant gradient descent." Proceedings of the 31st International Conference on Neural Information Processing Systems. 2017.

[3]. Sattler, Felix, et al. "On the byzantine robustness of clustered federated learning." ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2020.



## Experiments: Attack Schemes

- **Label Shuffling:** the labels of all samples are shuffled randomly in each corrupted client
- **Label Flipping:** the labels of all samples are switched to a random one in each corrupted client
- **Noisy Clients:** for CIFAR-10 and FEMNIST, first normalize the inputs into  $[0, 1]$ , then add Gaussian noise following  $N(0, 0.7)$  to the values, at last normalize them again; for Shakespeare, randomly select half of the characters and shuffle them so that the input sentence might be disordered
- Two corruption level: 30% or 50% are malicious clients.





## Experiments: Results

**Table 2. Averaged test accuracy over five random seeds for FedAvg, RFA, MKRUM, CFL and ARFL in four different scenarios**

CIFAR-10	Clean	Shuffling		Flipping		Noisy	
Corr. Per.	-	30%	50%	30%	50%	30%	50%
FedAvg	73.59 $\pm$ 0.44	61.17 $\pm$ 1.81	47.00 $\pm$ 7.51	65.01 $\pm$ 2.38	51.75 $\pm$ 7.75	<b>73.75 <math>\pm</math> 0.49</b>	73.61 $\pm$ 0.53
RFA	71.36 $\pm$ 0.47	57.86 $\pm$ 3.22	40.26 $\pm$ 9.14	55.47 $\pm$ 5.17	40.91 $\pm$ 11.06	73.74 $\pm$ 0.52	<b>73.69 <math>\pm</math> 0.63</b>
MKRUM	67.03 $\pm$ 0.93	59.27 $\pm$ 9.34	52.32 $\pm$ 14.90	60.21 $\pm$ 5.73	47.96 $\pm$ 10.25	73.41 $\pm$ 0.69	73.49 $\pm$ 0.49
CFL	71.68 $\pm$ 0.36	52.54 $\pm$ 1.71	50.29 $\pm$ 1.95	52.87 $\pm$ 1.07	51.67 $\pm$ 0.92	54.97 $\pm$ 1.14	55.26 $\pm$ 1.96
ARFL	73.42 $\pm$ 0.40	<b>71.68 <math>\pm</math> 1.01</b>	<b>69.66 <math>\pm</math> 0.73</b>	<b>71.78 <math>\pm</math> 0.53</b>	<b>70.25 <math>\pm</math> 0.56</b>	73.48 $\pm$ 0.56	73.29 $\pm$ 0.79

FEMNIST	Clean	Shuffling		Flipping		Noisy	
Corr. Per.	-	30%	50%	30%	50%	30%	50%
FedAvg	82.12 $\pm$ 0.20	61.91 $\pm$ 21.33	39.69 $\pm$ 20.80	70.19 $\pm$ 10.17	48.53 $\pm$ 23.49	79.94 $\pm$ 0.36	78.27 $\pm$ 0.47
RFA	82.11 $\pm$ 0.32	74.36 $\pm$ 7.52	52.02 $\pm$ 22.51	73.80 $\pm$ 7.49	50.75 $\pm$ 19.91	80.45 $\pm$ 0.30	79.21 $\pm$ 0.41
MKRUM	79.38 $\pm$ 0.41	57.51 $\pm$ 21.17	42.40 $\pm$ 24.84	78.57 $\pm$ 4.83	67.10 $\pm$ 7.35	<b>81.52 <math>\pm</math> 0.53</b>	<b>79.80 <math>\pm</math> 0.22</b>
CFL	82.18 $\pm$ 0.30	81.24 $\pm$ 0.47	36.03 $\pm$ 36.38	81.22 $\pm$ 0.36	65.54 $\pm$ 26.94	80.13 $\pm$ 0.70	79.21 $\pm$ 0.64
ARFL	<b>82.32 <math>\pm</math> 0.19</b>	<b>81.60 <math>\pm</math> 0.31</b>	<b>81.35 <math>\pm</math> 0.43</b>	<b>81.87 <math>\pm</math> 0.22</b>	<b>81.30 <math>\pm</math> 0.24</b>	80.71 $\pm$ 0.28	79.40 $\pm$ 0.45

Shakespeare	Clean	Shuffling		Flipping		Noisy	
Corr. Per.	-	30%	50%	30%	50%	30%	50%
FedAvg	53.80 $\pm$ 0.33	51.98 $\pm$ 0.48	47.70 $\pm$ 4.96	52.08 $\pm$ 0.39	41.85 $\pm$ 16.18	51.85 $\pm$ 0.56	50.43 $\pm$ 1.19
RFA	<b>54.27 <math>\pm</math> 0.41</b>	50.16 $\pm$ 1.28	32.49 $\pm$ 13.81	50.50 $\pm$ 1.02	23.84 $\pm$ 21.78	<b>52.17 <math>\pm</math> 0.50</b>	50.69 $\pm$ 1.04
MKRUM	50.81 $\pm$ 0.85	40.38 $\pm$ 7.44	24.46 $\pm$ 6.88	44.95 $\pm$ 2.43	16.11 $\pm$ 15.46	48.19 $\pm$ 0.40	45.67 $\pm$ 0.46
CFL	54.01 $\pm$ 0.34	49.76 $\pm$ 4.47	43.68 $\pm$ 12.68	51.09 $\pm$ 1.36	37.30 $\pm$ 19.76	51.98 $\pm$ 1.03	50.38 $\pm$ 1.39
ARFL	53.52 $\pm$ 0.32	<b>52.85 <math>\pm</math> 0.49</b>	<b>51.61 <math>\pm</math> 0.68</b>	<b>52.82 <math>\pm</math> 0.48</b>	<b>51.74 <math>\pm</math> 0.69</b>	52.09 $\pm$ 1.27	<b>50.98 <math>\pm</math> 0.75</b>

2.2

**Dynamic Quantifier  
*federated aggregation*  
operator (FL-IOWA-DQ)**



## Aggregation of FL-IOWA-DQ

$$\alpha_i = Q_{a,b,c,y_b}\left(\frac{i}{N}\right) - Q_{a,b,c,y_b}\left(\frac{i-1}{N}\right)$$

$$Q_{a,b,c,y_b}(x) = \begin{cases} 0 & 0 \leq x \leq a \\ \frac{x-a}{b-a} \cdot y_b & a \leq x \leq b \\ \frac{x-b}{c-b} \cdot (1-y_b) + y_b & b \leq x \leq c \\ 1 & c \leq x \leq 1 \end{cases}$$

where  $\alpha_i$  is corresponding to the update with **ith highest accuracy in validation dataset**.

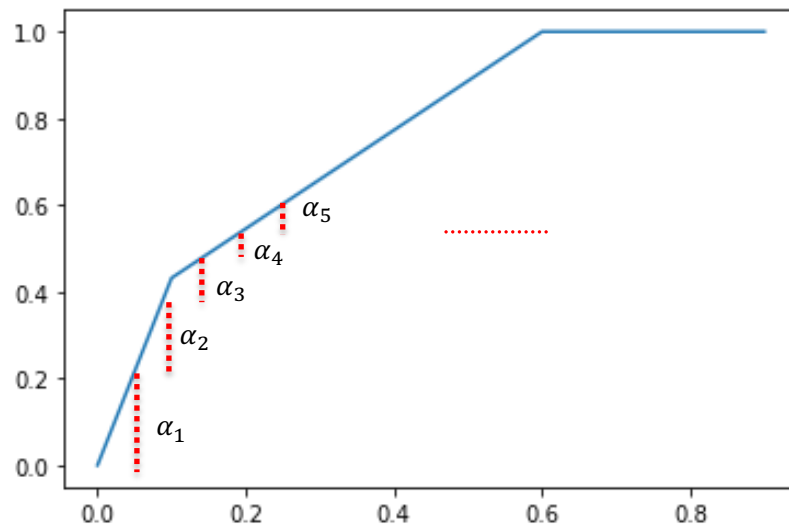
b and c are dynamic;  $b = 0.2c$ ;

c is the portion of clients whose distance to highest acc is least than  $\frac{3}{4}$  the maximum distance between lowest and highest one.

**Higher accuracy, higher weights.**



## Aggregation of FL-IOWA-DQ



**Fig. 1. An example of IOWA-DQ.  $a = 0$ ,  $b = 0.2c$ ,  $y_b = 0.4$ .**

$a=0$ : all clients with high acc are included in the aggregation;  
 $y_b$ : total weights assigned to updates with top acc;  
 $c$ : the portion of clients with non-zero weights



## Experiments: Models

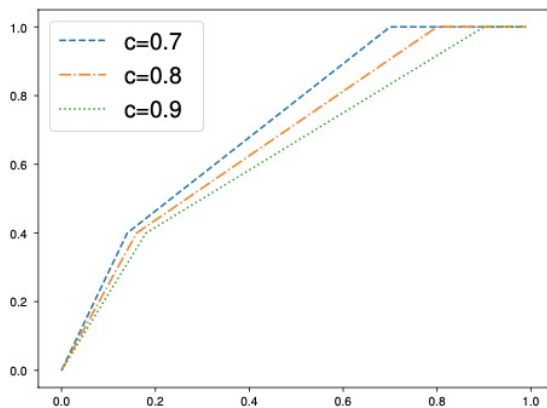
	a	b	c	y <sub>b</sub>
FL-IOWA-DQ-0.4	0	0.2*c	dynamic	0.4
FL-IOWA-DQ-0.75				0.75
IOWA-SQ-0.4		0.2	0.8	0.4
IOWA-SQ-0.75				0.75
FL-AL-80		0.8		1

**Table 3. Models with different configurations**

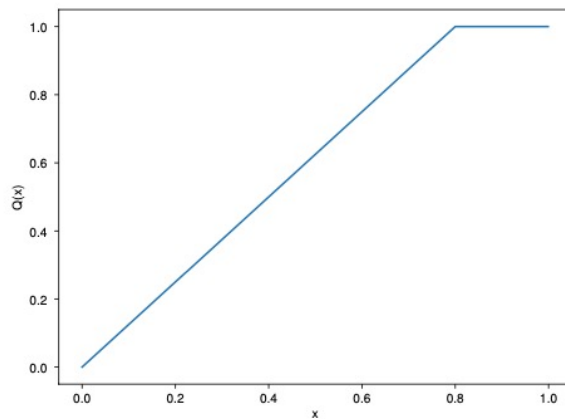
- Besides above, FedAvg and weighted FedAvg(W-FedAvg) are also used



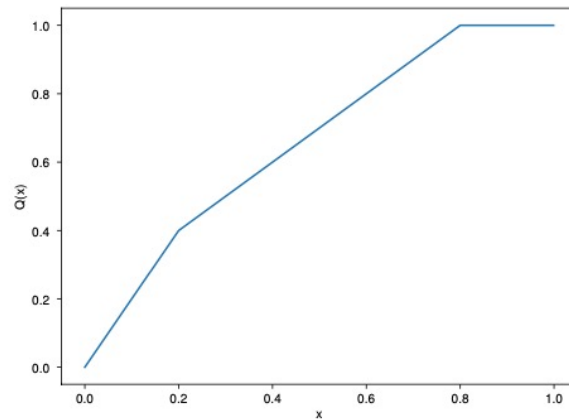
## Experiments: Models



**Fig. 2: FL-IOWA-DQ-0.4 with different  $c$ .**



**Fig. 3: FL-AL-80**



**Fig. 4: FL-IOWA-SQ-0.4**



## Experiments: Datasets

	EMNIST	Fashion MNIST
#Classes:	10;	
#Clients:	20/50;	
#Training:	200,000;	50,000;
#Validation:	40,000;	10,000
#Test:	40,000;	10,000
i.i.d:	No;	
Model:	CNN with two convolutional layers;	

- **Non i.i.d in size:** “randomly assign instances of a reduced number of labels to each client ”
- **Validation dataset:** follows the same distribution of the training subsets



## Experiments: Attack Scheme

*Definition 3.1 (Adversarial client):* Let  $C_i \in \{C_1, \dots, C_n\}$  be an arbitrary client of a FL environment whose original training dataset is  $D_i = \langle x_i^l; y_i^l \rangle$ , where  $x_i^l$  is the sample data and  $y_i^l$  the label. We say that  $C_i$  is an **adversarial client** if it uses the altered dataset  $D'_i$  as training dataset with

$$D'_i = \langle x_i^l; y_i^{\sigma(l)} \rangle,$$

where  $\sigma$  is a random permutation.

Attach Scheme: **Label Shuffling**





## Experiments: Results

1. AD Scenario: 10% of clients are malicious, 2 out of 20 or 5 out of 50.

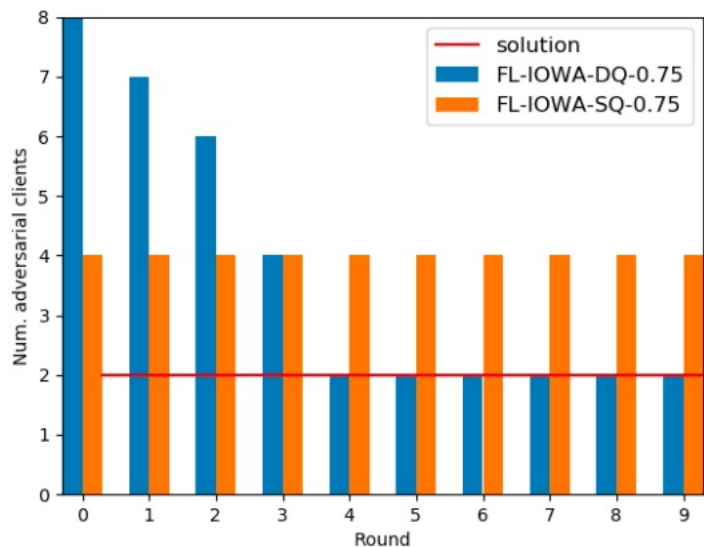
Table 4. Accuracy of Models in AD Scenario.

	EMNIST		Fashion-MNIST	
	20 clients	50 clients	20 clients	50 clients
<b>FL-FedAvg</b>	0.9826	0.9791	0.8661	0.8439
<b>FL-W-FedAvg</b>	0.9776	0.9774	0.8699	0.8321
<b>FL-AL-80</b>	0.9832	0.9803	0.8708	0.8469
<b>FL-IOWA-SQ-0.4</b>	0.9863	0.9824	0.8747	0.8541
<b>FL-IOWA-SQ-0.75</b>	0.9883	0.9869	0.8656	0.8671
<b>FL-IOWA-DQ-0.4</b>	0.9870	0.9886	<b>0.8782</b>	0.8694
<b>FL-IOWA-DQ-0.75</b>	<b>0.9900</b>	<b>0.9898</b>	0.8680	<b>0.8729</b>

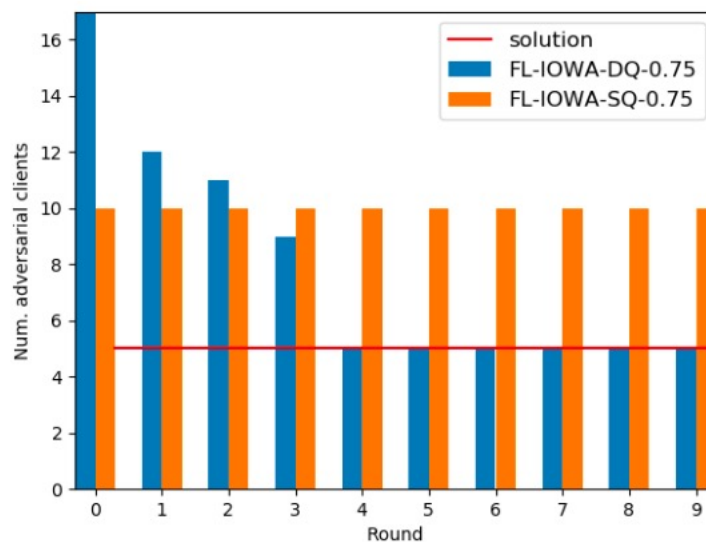


## Experiments: Results

1. AD Scenario: 10% of clients are malicious, 2 out of 20 or 5 out of 50.



**Fig. 5. Adversarial clients detected in AD Scenario with 20 clients**



**Fig. 6. Adversarial clients detected in AD Scenario with 50 clients**



## Experiments: Results

### 2. NON-AD Scenario: without adversarial clients

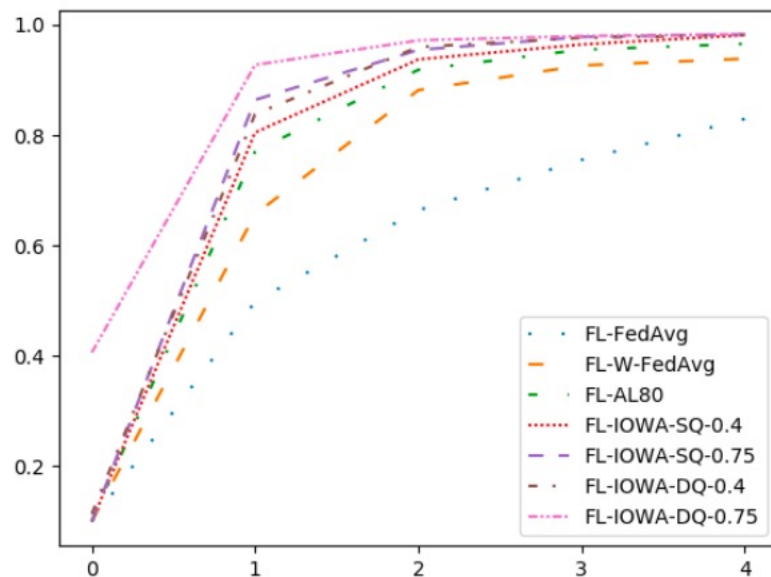
Table 5. Accuracy of Models in NON-AD Scenario.

	EMNIST		Fashion-MNIST	
	20 clients	50 clients	20 clients	50 clients
<b>FL-FedAvg</b>	0.9864	0.9801	0.8704	0.8452
<b>FL-W-FedAvg</b>	0.9857	0.9769	0.8721	0.8396
<b>FL-AL-80</b>	0.9861	0.9807	0.8772	0.8492
<b>FL-IOWA-SQ-0.4</b>	0.9882	0.9836	0.8793	0.8547
<b>FL-IOWA-SQ-0.75</b>	0.9890	0.9868	0.8726	0.8673
<b>FL-IOWA-DQ-0.4</b>	0.9891	0.9848	<b>0.8953</b>	0.8684
<b>FL-IOWA-DQ-0.75</b>	<b>0.9893</b>	<b>0.9873</b>	0.8923	<b>0.8728</b>



## Experiments: Results

### 2. NON-AD Scenario: without adversarial clients



**Fig. 7. Accuracy per round of FL models using 20 clients without adversarial clients (NON-AD Scenario) during the first 5 rounds**



## Experiments: Results

3. High-AD Scenario: 30% of clients are malicious, 6 out of 20 or 15 out of 50.

Table 6. Accuracy of Models in High-AD Scenario.

	EMNIST		Fashion-MNIST	
	20 clients	50 clients	20 clients	50 clients
<b>FL-FedAvg</b>	0.9788	0.9753	0.8451	0.8435
<b>FL-W-FedAvg</b>	0.9769	0.9758	0.8456	0.8228
<b>FL-AL-80</b>	0.9713	0.9781	0.8439	0.8212
<b>FL-IOWA-SQ-0.4</b>	0.9826	0.9820	0.8468	0.8539
<b>FL-IOWA-SQ-0.75</b>	0.9844	0.9861	0.8518	0.8604
<b>FL-IOWA-DQ-0.4</b>	<b>0.9876</b>	0.9860	0.8648	0.8610
<b>FL-IOWA-DQ-0.75</b>	0.9873	<b>0.9874</b>	<b>0.8722</b>	<b>0.8684</b>

---

3

## **Extra Experiments**

---



## Setting

<b>Datasets</b>	Fashion MNIST(fMNIST) <ul style="list-style-type: none"><li>• 10 classes</li><li>• 60k/10k</li></ul>	CIFAR-10 <ul style="list-style-type: none"><li>• 10 classes</li><li>• 50k/10k</li></ul>
<b>Date Distribution</b>	I.I.D. <ul style="list-style-type: none"><li>• # local epochs = 5</li></ul>	Extreme NON-I.I.D[1] <ul style="list-style-type: none"><li>• # local epochs = 1</li><li>• #classes per client = 2</li></ul>
<b>Attack schemes</b>	Label Shuffling <ul style="list-style-type: none"><li>• 10% corrupted data</li></ul>	<b>Label Mislabeling*</b> <ul style="list-style-type: none"><li>• One class</li></ul>
<b>NN</b>	CNN with two convolutional layers	
<b>#max aggregation rounds</b>	40	
<b>#clients</b>	10	
<b>#mal clients</b>	1	

\***Label Mislabeling**: in mal client, randomly select an original class and replace its labels with target class



## Setting: Extreme Non-i.i.d

Data split:

```
- Client 0: [ 857 1500 857 1000 1000 1200 3000 0 1000 2000]
- Client 1: [ 857 1500 857 1000 1000 1200 0 0 1000 2000]
- Client 2: [ 857 0 857 0 0 0 0 0 1000 0]
- Client 3: [ 0 0 857 0 0 0 0 0 0 0]
- Client 4: [ 857 0 0 1000 0 0 0 0 0 0]
- Client 5: [ 857 0 0 1000 1000 0 0 0 0 0]
- Client 6: [ 857 0 857 0 1000 1200 0 0 1000 0]
- Client 7: [ 0 1500 857 1000 1000 1200 0 0 1000 0]
- Client 8: [ 858 1500 858 1000 1000 1200 0 0 1000 2000]
- Client 9: [ 0 0 0 0 0 0 3000 6000 0 0]
- Data size for clients: [12414, 9414, 2714, 857, 1857, 2857, 4914, 6557, 9416, 9000]
```

VS

Data split:

```
- Client 0: [ 1 0 0 0 0 0 0 0 2614 2614]
- Client 1: [ 0 0 0 0 0 0 3165 3165 0 0]
- Client 2: [2201 1 0 0 0 0 0 0 0 2201]
- Client 3: [ 0 0 0 0 4449 4449 1 0 0 0]
- Client 4: [ 0 0 0 2706 1551 1155 0 0 0 0]
- Client 5: [ 0 0 0 0 0 396 2834 2835 1366 0]
- Client 6: [ 604 0 0 0 0 0 0 0 1788 1185]
- Client 7: [ 0 2247 2247 1 0 0 0 0 0 0]
- Client 8: [3194 3752 1904 0 0 0 0 0 232 0]
- Client 9: [ 0 0 1849 3293 0 0 0 0 0 0]
- Data size for clients: [5229, 6330, 4403, 8899, 5412, 7431, 3577, 4495, 9082, 5142]
```

**Fig. 8. Data split of fMNIST in non-i.i.d in size manner.**

**Fig. 9. Data split of fMNIST in extreme non-i.i.d manner. Each client has no more than four classes.**

**Extreme Non-i.i.d:**

- Has less variance of the size of the local samples;
- Has sparser distribution of classes;

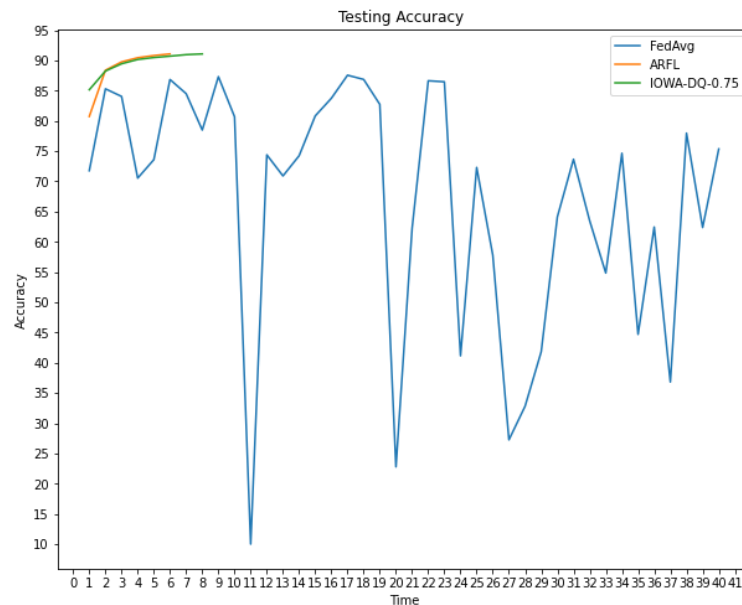




## Results: fMNIST-iid



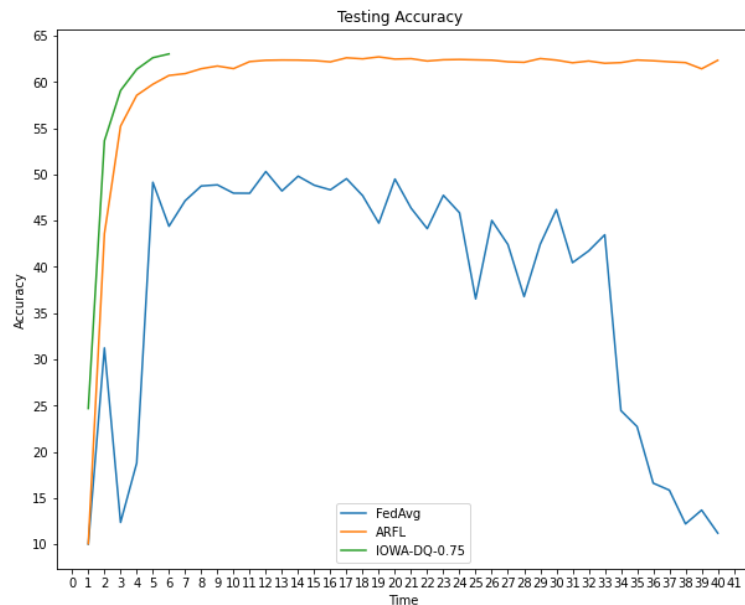
**Fig. 10. Testing accuracy of models in shuffling scenario.**



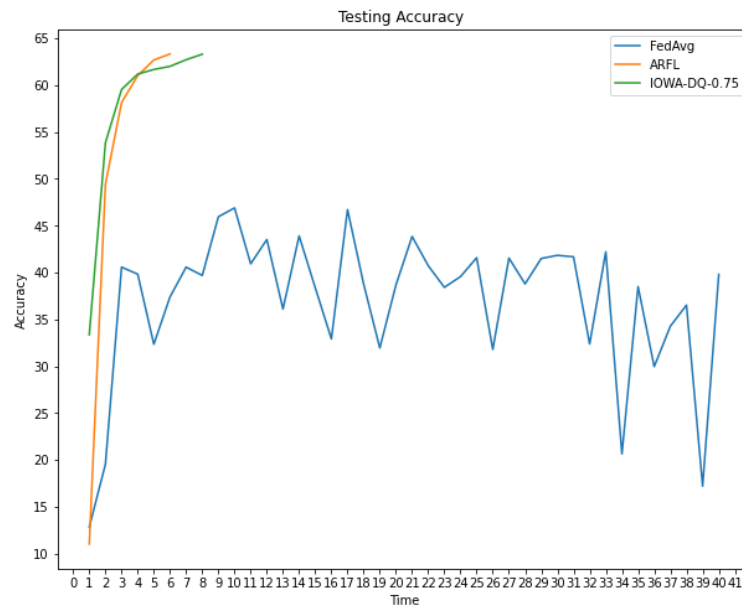
**Fig. 11. Testing accuracy of models in mislabeling scenario.**



## Results: CIFAR-10-iid



**Fig. 12. Testing accuracy of models in shuffling scenario.**



**Fig. 13. Testing accuracy of models in mislabeling scenario.**



## Results: fMNIST-non-iid



**Fig. 14. Testing accuracy of models in shuffling scenario.**



**Fig. 15. Testing accuracy of models in mislabeling scenario.**



## Results: CIFAR-non-iid



**Fig. 16. Testing accuracy of models in shuffling scenario.**



**Fig. 17. Testing accuracy of models in mislabeling scenario.**



## Conclusions

---

- In iid setting, IOWA-DQ and ARFL are both effective and protect the model from serious performance degeneration. They have very close performance.
- In extreme non-iid setting, defenses by simply adjusting weights are ineffective, even weaker than FedAvg.
- The data distribution, corruption degree highly affect the robustness of defenses.



## References

- Li, Shenghui, et al. "Auto-weighted Robust Federated Learning with Corrupted Data Sources." *arXiv preprint arXiv:2101.05880* (2021).
- Rodríguez-Barroso, Nuria, et al. "Dynamic federated learning model for identifying adversarial clients." *arXiv preprint arXiv:2007.15030* (2020).
- Kairouz, Peter, et al. "Advances and open problems in federated learning." *arXiv preprint arXiv:1912.04977* (2019).
- Learning Multiple Layers of Features from Tiny Images, Alex Krizhevsky, 2009.
- Caldas, Sebastian, et al. "Leaf: A benchmark for federated settings." *arXiv preprint arXiv:1812.01097* (2018).
- Pillutla, Krishna, Sham M. Kakade, and Zaid Harchaoui. "Robust aggregation for federated learning." *arXiv preprint arXiv:1912.13445* (2019).
- Blanchard, Peva, et al. "Machine learning with adversaries: Byzantine tolerant gradient descent." *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 2017.
- Sattler, Felix, et al. "On the byzantine robustness of clustered federated learning." *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020.
- Zhao, Yue, et al. "Federated learning with non-iid data." *arXiv preprint arXiv:1806.00582* (2018).



**Thanks**