

# FL-WBC: Enhancing Robustness against Model Poisoning Attacks in Federated Learning from a Client Perspective

Yuri Dimitre Dias de Faria

y218172@dac.unicamp.br

18 de Junho 2024

# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# Table of Contents

1. **Introdução**
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# Introdução

## Funcionamento do Aprendizado Federado

- Definição
  - Abordagem de aprendizado de máquina distribuída.
  - Permite a colaboração de múltiplos dispositivos de borda no treinamento de um modelo global.
- Funcionamento
  - Treinamento Local: Cada dispositivo treina um modelo localmente usando seus próprios dados.
  - Envio de Atualizações: Dispositivos enviam apenas atualizações de modelo (gradientes ou parâmetros) para um servidor central.
  - Agregação Centralizada: O servidor central agrega as atualizações recebidas para criar um modelo global atualizado.
  - Distribuição do Modelo: O modelo global atualizado é enviado de volta aos dispositivos.
  - Ciclo Iterativo: O processo de treinamento local, envio de atualizações e agregação é repetido várias vezes.

# Introdução

## Problemas de segurança em FL

- Alta vulnerabilidade a ataques de envenenamento (*poisoning attacks*) em sistema FL.

# Introdução

## Objetivo

- Melhorar a robustez contra esse tipo de ataques por mecanismo de defesa *client-based*.

# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# Background

## Model Poisoning Attacks

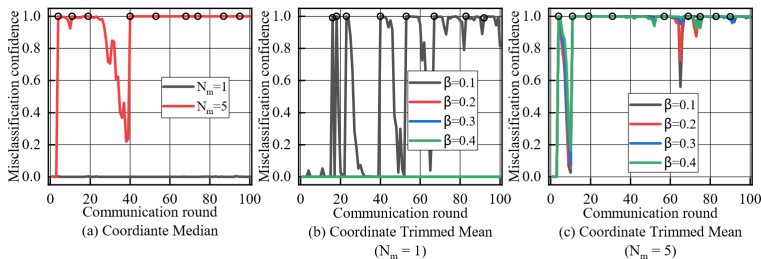
Tipos de ataques em modelos em FL.

- Não-direcionado (untargeted): Ataca o modelo global indiscriminadamente para que ele tenha uma alta taxa de erro.
- Direcionado (targeted): Ataca o modelo global para gerar classificações incorretas desejadas pelo invasor.



# Background

## Defesas da literatura



**Figure:** Desempenho de *robust aggregation* para ataques em grande escala

# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# Contribuições

## Attack Effect on Parameter (AEP)

AEP  $\delta_t$  é uma métrica do acúmulo das mudanças do modelo global até a rodada  $t$  devido a ataques conduzidos.

$$W_t(\mathbb{S}_i/\mathbb{M}) \leftarrow \frac{N}{K} \sum_{k \in \mathbb{S}_t} p^k W_{t,l}^k(\alpha = 1)$$

$$\delta_t \triangleq W_t(\mathbb{S}/M) - W_t$$

Onde  $\mathbb{M}$  é um conjunto de atacantes,  $W_t(\mathbb{S}_i/\mathbb{M})$  representa os pesos do modelo global na rodada  $t$  quando todos os dispositivos maliciosos em  $\mathbb{S}_i (i \leq t)$  não performaram um ataque em  $i$  rodadas,  $N$  representa o número de dispositivos e  $K$  é o número de dispositivos selecionados para agregação na rodada.

# Contribuições

## *Attack Effect on Parameter (AEP)*

- Por que pode ser difícil eliminar um envenenamento do modelo global?
- Razão: A transmissão do AEP  $\delta_t$  para o modelo global é determinada por  $H_{t,i}^k$  (matriz hessiana), inacessível pelo servidor central.

# Contribuições

## Mecanismo de defesa FL-WBC

- Proposto um mecanismo de perturbação em matrizes hessianas em dispositivos benignos para combater um envenenamento,

# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# FL-WBC

## Objetivos

- Ideia central: Perturbar o espaço de parâmetros no treinamento local para mitigar efeitos de envenenamento.
- Objetivo 1: Para manter a desempenho da tarefa benigna, o *loss* da mesma deve ser minimizado.
- Objetivo 2: Para prevenir que AEP se esconda no *kernel* das matrizes hessianas em dispositivos benignos, o *kernel*  $H_{l+1,i}^k$  deve ser perturbado.

Para perturbar  $H_{l,i}^k$ , foi considerado a segunda derivada parcial da função de perda, onde a diagonal descreve a mudança do gradiente pelas iterações.

# FL-WBC

## Formulação matemática

$$W_{t,i+1}^{\hat{k}} = W_{t,i}^k - \eta_{t,i} \nabla F^k(W_{t,i}^k, \xi_{t,i}^k)$$

$$W_{t,i+1}^k = W_{t,i}^{\hat{k}} + \eta_{t,i} \Upsilon_{t,i}^k \odot M_{t,i}^k$$

Onde  $\Upsilon_{t,i,r,c}^k$  é uma matriz ruído de mesma dimensão de  $W$ ,  $M_{t,i}^k$  é uma máscara binária, na qual os elementos são determinados por:

$$M_{t,i,r,c}^k = \begin{cases} 1, & |(W_{t,i+1}^{\hat{k}} - W_{t,i}^k) - \Delta W_{t,i}^k|_{r,c} / \eta_{t,i} \leq |\Upsilon_{t,i,r,c}^k| \\ 0, & |(W_{t,i+1}^{\hat{k}} - W_{t,i}^k) - \Delta W_{t,i}^k|_{r,c} / \eta_{t,i} > |\Upsilon_{t,i,r,c}^k| \end{cases}$$

Nesse trabalho, foi escolhido  $\Upsilon_{t,i,r,c}^k$  como sendo um ruído laplaciano com média zero e desvio padrão  $s$ .



---

**Algorithm 1** Local training process applying FL-WBC on a benign device in round  $t$ .

---

**Input:** Local training data  $\mathbb{D} \in \mathbb{R}^{L \times P \times Q}$ ; Local objective function  $F : \mathbb{R}^{P \times Q} \rightarrow \mathbb{R}$ ; Local model parameters  $\mathbf{W} \in \mathbb{R}^{M \times N}$ ; The number of local training iterations  $I$ ; Learning rates  $\eta_{t,i}$  for  $i \in [I]$ ; Standard deviation of Laplace noise  $s$ .

**Output:** Learnt model parameter  $\mathbf{W}$  with our defense.

- 1: Initialize  $\mathbf{W}_{-1}, \mathbf{W}_{-2}$ ;
- 2:  $i \leftarrow 0$ ;
- 3: **for** batch  $\mathcal{B}$  in  $\mathbb{D}$  **do**
- 4:     Randomly generate a Laplace noise matrix  $\Upsilon \in \mathbb{R}^{M \times N}$  with  $mean = 0$  and  $std = s$ ;
- 5:      $\mathbf{W}_{-1} \leftarrow \mathbf{W}$ ;
- 6:      $\mathbf{W} \leftarrow \mathbf{W} - \eta_{t,i} \nabla F(\mathbf{W}, \mathcal{B})$ ;
- 7:     **if** this is not the first training batch **then**
- 8:          $\mathbf{W}^* \leftarrow (\mathbf{W} - \mathbf{W}_{-1}) - (\mathbf{W}_{-1} - \mathbf{W}_{-2})$ ;
- 9:         Find the set  $\mathbb{S}$  which contains the indices of elements in  $|\mathbf{W}^*| - \eta_{t,i} |\Upsilon|$  which are less than or equal to 0;
- 10:         **for**  $j, k \in \mathbb{S}$  **do**
- 11:              $\mathbf{W}_{j,k} \leftarrow \mathbf{W}_{j,k} + \eta_{t,i} \Upsilon_{j,k}$ ;
- 12:         **end for**
- 13:     **end if**
- 14:      $\mathbf{W}_{-2} \leftarrow \mathbf{W}_{-1}$ ;
- 15:      $i \leftarrow i + 1$
- 16: **end for**

---

Figure: Algoritmo FL-WBC

# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

# Procedimentos Experimentais

## Configurações e datasets

- Datasets: FashionMNIST e CIFAR10.
- Configurações IID e non-IID.
- Dataset  $D_m$  malicioso com samples de uma e múltiplas imagens.
- Comparação com métodos de agregação robusta (CMA e CTMA) e *differential privacy* (DP).
- DP: *Central Differential Privacy* (CDP) e *Local Differential Privacy* (LPD).
- Rodadas contam com 10 dispositivos, onde 5 deles são maliciosos caso a rodada for maligna.
- Probabilidade de 10% de ser uma rodada adversária.

# Procedimentos Experimentais

## Resultados Robust Aggregation com 1 imagem

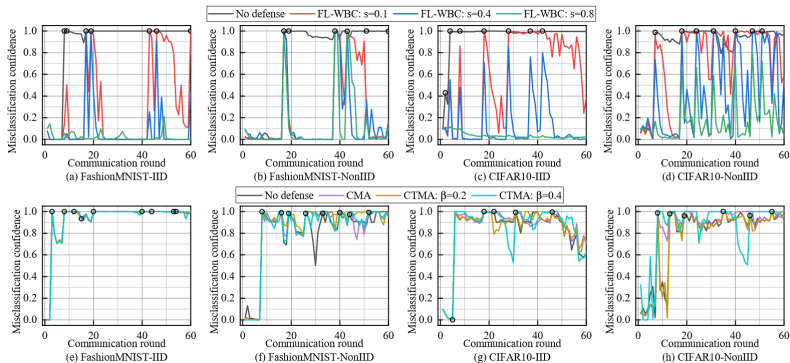


Figure: Comparação da confiança na classificação incorreta

# Procedimentos Experimentais

## Resultados DP 1 imagem

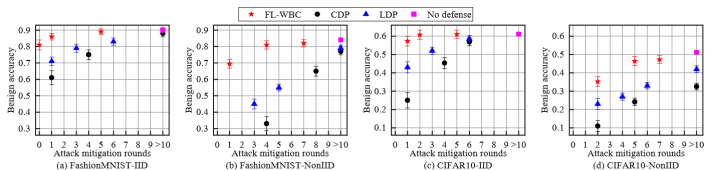


Figure: Comparação da acurácia benigna e rodadas de mitigação de ataques

# Procedimentos Experimentais

## Resultados Robust Aggregation multiplas imagem

Defense	Fashion-MNIST (IID)	Fashion-MNIST (non-IID)	CIFAR10 (IID)	CIFAR10 (non-IID)
CTMA ( $\beta = 0.1$ )	7	9	8	>10
CTMA ( $\beta = 0.2$ )	7	8	8	9
CTMA ( $\beta = 0.4$ )	6	8	7	9
CMA	5	7	6	8

Figure: Comparação de rodadas de mitigação de ataques

# Procedimentos Experimentais

## Resultados DP multiplas imagem

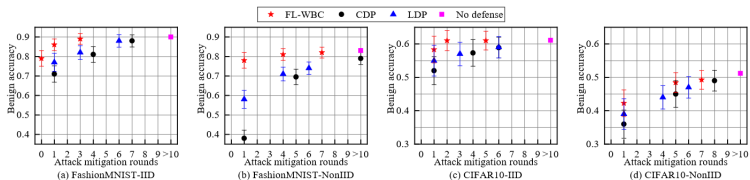
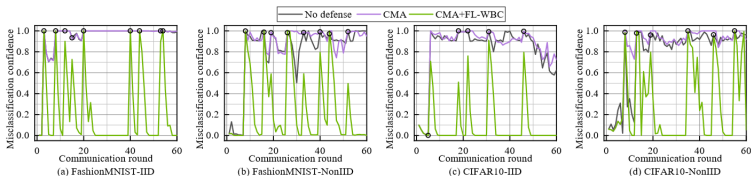


Figure: Comparação da acurácia benigna e rodadas de mitigação de ataques

# Procedimentos Experimentais

## Integração do FL-WBC com Robust Aggregation



**Figure:** Comparação da confiança da classificação incorreta e rodadas de comunicação



# Table of Contents

1. Introdução
2. Background
3. Contribuições
4. FL-WBC
5. Procedimentos Experimentais
6. Conclusão

## Conclusões

- Defesa contra ataque de envenenamento de modelo baseada no cliente .
- Resultados demonstram que o sistema supera as linhas de base de mitigação do efeito do ataque.
- A defesa consegue proteger contra o ataque em menos rodadas de comunicação que outros tipos de defesas com uma menor degradação da utilidade do modelo.
- Pode ser estendido para outros tipos de ataques de envenenamento, como backdoor.