



# Rethinking the Flat Minima Searching in Federated Learning

Taehwan Lee<sup>1</sup>

Sung Whan Yoon<sup>1,2</sup>

- Published as a conference paper at ICML 2024

**汇报人：游先耀**

**2024.10.15**

- **平滑最小Loss平面搜索：**将loss最小值附近的平坦度整合到成本函数中，有助于在训练期间找到更平坦的局部模型。该方法已被证实可以增强联邦学习算法的性能，尤其是在异构环境中。

**SAM optimizer:** The SAM optimizer transforms a loss function  $f(w)$  into a min-max cost function as follows:

$$\min_w \max_{\|\delta\| \leq \rho} F(w + \delta), \quad (4)$$

where  $\rho$  is a positive real number and  $\|\delta\|$  is L2-norm of  $\delta$ .

- **FedSAM:** 该方法将 SAM 优化器直接应用于联邦学习中的本地训练，旨在通过优化局部目标来找到更扁平的局部模型，当局部模型聚合到全局模型中时，有助于提高性能。

**FedSAM:** By adopting the min-max problem of Eq. (4) in local training, FedSAM perturbs local model  $w_{i,k}^r$ :

$$\tilde{w}_{i,k}^r = w_{i,k}^r + \delta = w_{i,k}^r + \rho g_{i,k}^r / \|g_{i,k}^r\| \quad (5)$$

$$w_{i,k+1}^r = w_{i,k}^r - \eta l \tilde{g}_{i,k}^r, \quad (6)$$

- **平坦度差异:** 本文指出了一个问题称为平坦度差异的关键问题，该问题是指观察到在局部训练期间搜索最小平面并不能保证全局模型也将位于损失格局的平坦区域。这种差异可能导致聚合的全局模型的性能不理想，即使局部模型经过有效训练。

$$\Delta_{\mathcal{F}} := \left| \max_{\|\delta\| \leq \rho} F(w + \delta) - F(w) - \left[ \sum_{i=1}^N \frac{m_i}{m} \max_{\|\delta_i\| \leq \rho} F_i(w_i + \delta_i) - F_i(w_i) \right] \right|.$$

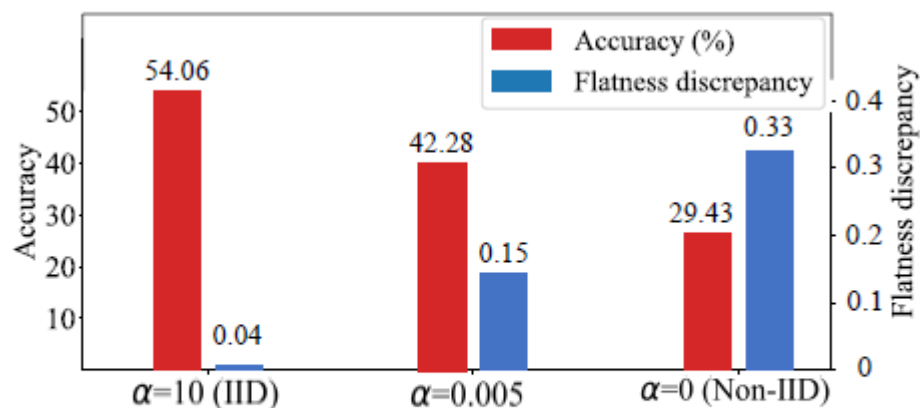


Figure 1: The performance and the flatness discrepancy ( $\Delta_{\mathcal{F}}$ ) of FedSAM for the CIFAR-100 experiment

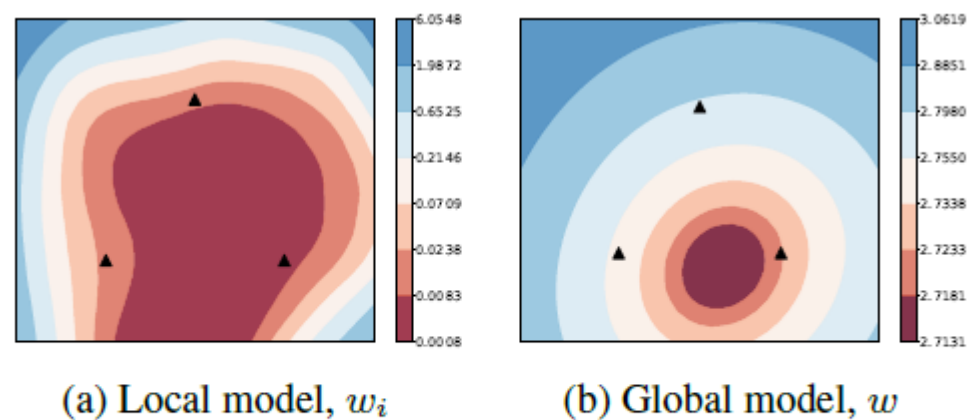


Figure 2: Visualization of the loss surface of FedSAM for the CIFAR-100 case ( $\alpha = 0$ ).

- **FedGF**: 为了解决这些问题, 本文提出了一种名为全球平坦度联邦学习 (FedGF) 的新方法。该方法旨在明确地为全局模型追求更平坦的最小值, 从而缓解平坦度差异并在异构 FL 基准测试中实现显著的性能改进。

$$g_{i,k}^r = \nabla F_i(w_{i,k}^r, \zeta_{i,k}) \quad (8)$$

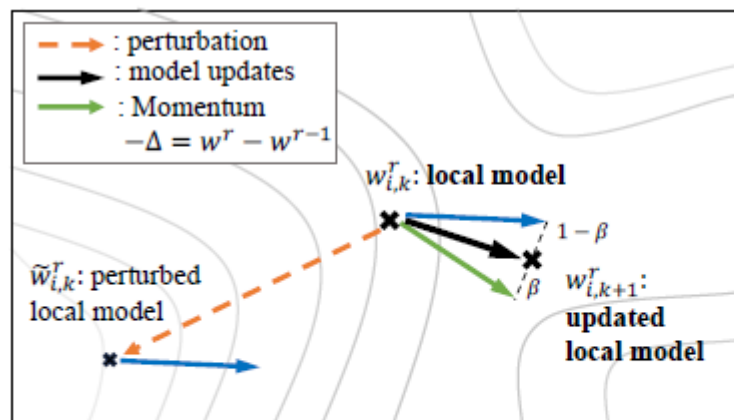
$$\tilde{w}_{i,k}^r = w_{i,k}^r + \rho g_{i,k}^r / \|g_{i,k}^r\| \quad (\text{perturbed local model}) \quad (9)$$

$$\Delta^r = w^{r-1} - w^r \quad (10)$$

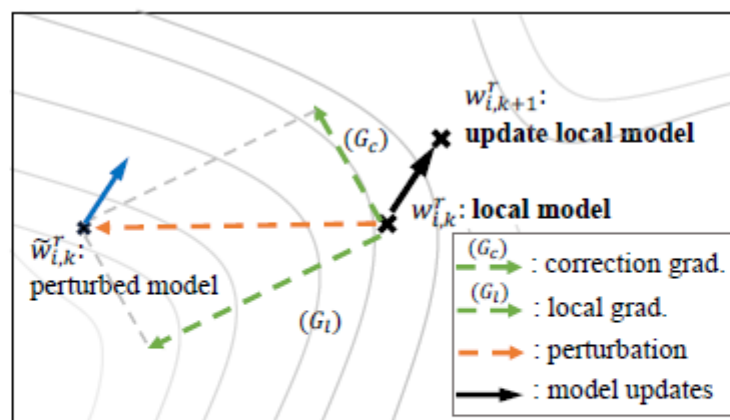
$$\tilde{w}^r = w^r + \rho \Delta^r / \|\Delta^r\| \quad (\text{perturbed global model}) \quad (11)$$

$$\tilde{w}_{i,k,c}^r = c\tilde{w}^r + (1-c)\tilde{w}_{i,k}^r, \quad (12)$$

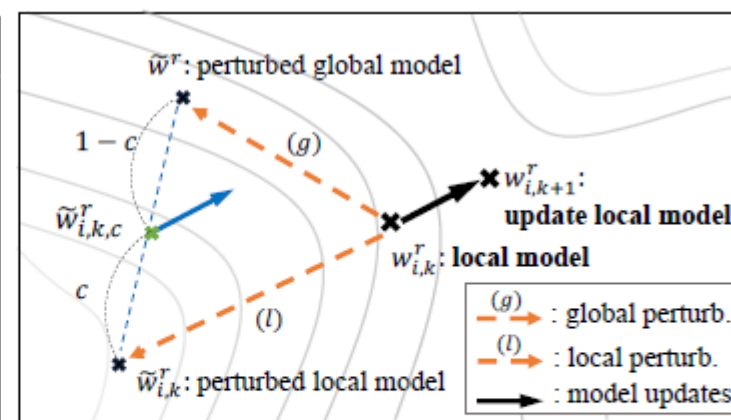
$$w_{i,k+1}^r = w_{i,k}^r - \eta_l \nabla F_i(\tilde{w}_{i,k,c}^r, \zeta_{i,k}). \quad (13)$$



(a) MoFedSAM (Qu et al., 2022)



(b) FedSMOO (Sun et al., 2023a)



(c) FedGF (ours)

Figure 3: Schematic of MoFedSAM, FedSMOO, and FedGF. The gray line illustrates the loss landscape of local distribution.

- **权重 $c$ 的选择:** 权重 $c$  控制局部和全局扰动之间的扰动模型，基于以下策略：当局部模型与全局模型偏差很大时（non-IID），为了确保全局模型的平坦度，使用较大的  $c$ ；否则（IID），使用较小的  $c$ 。

$$D^r = \frac{1}{|S^r|} \sum_{i \in S^r} \|w^r - w_{i,K}^r\|_2. \quad (15)$$

$$c = \frac{1}{W} \sum_{i=r-W+1}^r I^i. \quad (16)$$

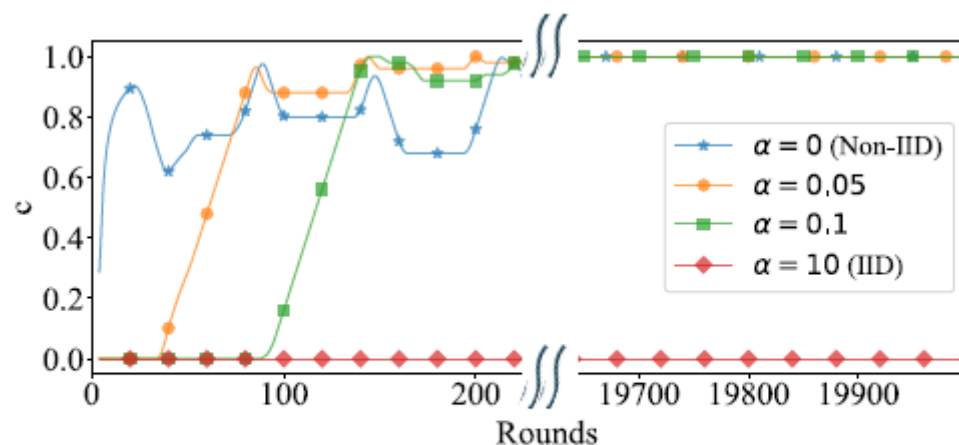


Figure 7: Behavior of  $c$  for CIFAR-100

# 实验：BENCHMARK RESULTS

Table 1: Test accuracies with of the FL algorithms on the CIFAR-10 and CIFAR-100 benchmarks

Task	Algorithms	Dirichlet distribution parameter $\alpha$								
		$Dir.(\alpha = 0, \text{non-IID})$			$Dir.(\alpha = 0.005)$			$Dir.(\alpha = 10, \text{IID})$		
		Number of participating clients per each round								
		5	10	20	5	10	20	5	10	20
CIFAR-10	FedAvg	63.63	65.83	68.33	67.85	71.37	73.03	82.90	82.96	82.93
	FedAvgM	62.73	65.61	68.57	67.56	71.32	75.53	82.72	83.60	83.30
	FedProx	63.13	65.95	67.98	68.06	71.42	72.87	82.72	83.19	82.92
	SCAFFOLD	(X)	(X)	(X)	57.13	56.46	45.27	82.93	83.05	83.39
	FedDyn	66.84	71.01	69.45	70.74	73.78	75.43	83.07	83.58	83.67
	FedSAM	68.11	71.17	72.49	71.87	74.31	76.07	83.78	83.88	83.82
	FedASAM	73.32	74.5	75.49	74.96	75.59	76.57	83.11	83.28	82.89
	MoFedSAM	73.1	71.08	76.66	74.43	77.53	79.27	80.9	81.01	81.02
	FedGAMMA	45.32	47.55	35.07	46.99	48.44	35.58	74.99	66.12	54.85
	FedSMOO	68.82	71.59	72.48	71.9	74.46	75.44	83.72	83.67	83.79
	<b>FedGF</b>	<b>78.41</b>	<b>79.68</b>	<b>80.86</b>	<b>78.79</b>	<b>79.39</b>	<b>79.69</b>	<b>84.71</b>	<b>83.94</b>	<b>83.85</b>
CIFAR-100	FedAvg	29.35	33.79	36.62	38.15	40.58	41.27	50.41	50.20	49.98
	FedAvgM	29.94	30.07	39.35	38.64	40.72	<b>48.44</b>	50.37	51.2	50.57
	FedProx	29.19	33.16	36.41	38.54	40.52	40.77	50.10	49.98	49.96
	SCAFFOLD	(X)	(X)	(X)	36.25	(X)	(X)	52.28	52.12	52.48
	FedDyn	(X)	(X)	(X)	(X)	(X)	(X)	51.74	52.41	52.59
	FedSAM	29.43	34.32	36.88	42.28	44.57	45.18	54.06	53.75	53.5
	FedASAM	34.43	37.09	38.93	44.36	45.76	46.94	<b>54.6</b>	54.42	<b>54.73</b>
	MoFedSAM	29.02	35.82	41.26	34.64	42.24	44.92	52.13	52.21	52.07
	FedGAMMA	(X)	(X)	(X)	20.52	14.76	10.33	47.43	38.18	25.06
	FedSMOO	35.35	38.78	40.82	44.39	46.03	47.5	54.31	54.89	54.65
	<b>FedGF</b>	<b>45.37</b>	<b>46.86</b>	<b>47.77</b>	<b>46.48</b>	<b>46.70</b>	46.08	54.16	<b>54.62</b>	54.59

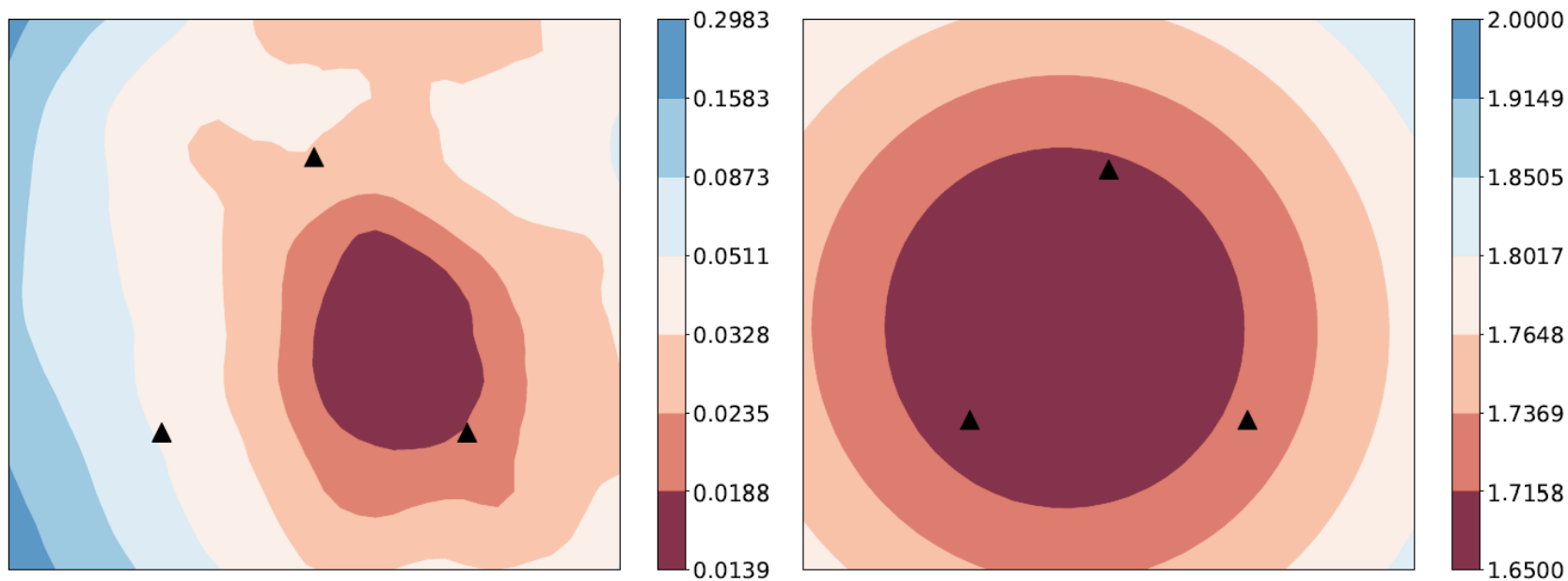
(X) indicates that the method fails to train, so the results remain at the same level as the random prediction.



# 实验：固定权重 $c$ vs 动态计算权重 $c$

Table 5: Static  $c$  vs. FedGF (adaptive  $c$ )

Dataset	IIDness	$c = 0$	$c = 0.5$	$c = 1$	<b>FedGF</b>
CIFAR-10	Non-IID	68.11	71.24	78.05	<b>78.41</b>
	IID	83.78	82.95	81.94	<b>84.71</b>
CIFAR-100	Non-IID	29.43	26.64	44.39	<b>45.37</b>
	IID	54.06	52.47	46.68	<b>54.16</b>



(a) Local model,  $w_i$

(b) Global model,  $w$

Figure 6: Loss surface of FedGF for CIFAR-100 ( $\alpha = 0$ ).