Improving Generalization in Federated Learning by Seeking Flat Minima

Debora Caldarola^{*1}, Barbara Caputo^{1,2}, and Marco Ciccone^{*1}

¹Politecnico di Torino, ²CINI name.surname@polito.it

Abstract. Models trained in federated settings often suffer from degraded performances and fail at generalizing, especially when facing heterogeneous scenarios. In this work, we investigate such behavior through the lens of geometry of the loss and Hessian eigenspectrum, linking the model's lack of generalization capacity to the sharpness of the solution. Motivated by prior studies connecting the sharpness of the loss surface and the generalization gap, we show that i) training clients locally with Sharpness-Aware Minimization (SAM) or its adaptive version (ASAM) and ii) averaging stochastic weights (SWA) on the server-side can substantially improve generalization in Federated Learning and help bridging the gap with centralized models. By seeking parameters in neighborhoods having uniform low loss, the model converges towards flatter minima and its generalization significantly improves in both homogeneous and heterogeneous scenarios. Empirical results demonstrate the effectiveness of those optimizers across a variety of benchmark vision datasets (e.g. CIFAR10/100, Landmarks-User-160k, IDDA) and tasks (large scale classification, semantic segmentation, domain generalization).

1 Introduction

Federated Learning (FL) [51] is a machine learning framework enabling the training of a prediction model across distributed clients while maintaining their privacy, never disclosing local data. In recent years it has had a notable resonance in the world of computer vision, with applications ranging from large-scale classification [27] to medical imaging [22] to domain generalization [49] and many others [43,72,21,66]. The learning paradigm is based on communication rounds where a sub-sample of clients trains the global model independently on their local datasets, and the produced updates are later aggregated on the server-side. The heterogeneous distribution of clients' data, which is usually non-i.i.d. and unbalanced, poses a major challenge in realistic federated scenarios, leading to degraded convergence performances [73,26,45]. Locally, the model has only access to a small portion of the data failing to generalize to the rest of the underlying distribution. That contrasts with the standard centralized training, where the

^{*} Equal contribution

Official code: https://github.com/debcaldarola/fedsam



Fig. 1: Cross-entropy loss landscapes of the global model in heterogeneous ($\alpha = 0$) and homogeneous ($\alpha = 1k$) federated scenarios on CIFAR100. When trained with FedAvg, the global model converges towards sharp minima. The sharpness-aware optimizer ASAM significantly smooths the surfaces.

learner can uniformly sample from the whole distribution. While many promising works in the literature focus on regularizing the local objective to align the global and local solutions, thus reducing the client drift [45,34,1], less attention has been given to the explicit optimization of the loss function for finding better minima. Several works studied the connection between the sharpness of the loss surface and model's generalization [25,35,41,38,59,30,12], and proposed effective solutions based on the minimization of the derived generalization bound [69,18,40] or on averaging the network's parameters along the trajectory of SGD [29].

In this work, we first analyze the heterogeneous federated scenario to highlight the causes behind the poor generalization of the federated algorithms. We hypothesize during local training the model overfits the current distribution, and the resulting average of the updates is strayed apart from local minima. Thus, the global model is not able to generalize to the overall underlying distribution and has a much slower convergence rate, *i.e.* it needs a much larger number of rounds to reach the performance of the homogeneous setting. To speed up training and reduce the performance gap in the case of non-i.i.d. data, we look at improving the generalization ability of the model. Motivated by recent findings relating the geometry of the loss and the generalization gap [35,16,41,32] and by the achievements in the field of Vision Transformers [12], we analyze the loss landscape in the federated scenario and find out that models converge to sharp minima (Fig.1), hence the poor generalization. As a solution, we introduce methods of the current literature that explicitly look for flat minima: i) Sharpness-Aware Minimization (SAM) [18] and its adaptive version (ASAM) [40] on the client-side and ii) Stochastic Weight Averaging (SWA) [29] on the server-side. These modifications, albeit simple, surprisingly lead to significant improvements. Their use is already effective if taken individually, but the best performance is obtained when combined. The resultant models exhibit smoother loss surfaces and improved final performance consistently across several vision tasks. To summarize, our main contributions are:

 We analyze the behavior of models trained in heterogeneous and homogeneous federated scenarios by looking at their convergence points, loss surfaces and Hessian eigenvalues, linking the lack in generalization to sharp minima. Improving Generalization in Federated Learning by Seeking Flat Minima

- To encourage convergence towards flatter minima, we introduce SAM and ASAM in the local client-side training and SWA in the aggregation of the updates on the server-side. The resultant models show smoother loss landscapes and lower Hessian eigenvalues, with improved generalization capacities.
- We test our approach on multiple vision tasks, *i.e.* small and large scale classification [27], domain generalization [7] and semantic segmentation [50,11].
- We compare our method with strong data augmentations techniques and state-of-the-art FL algorithms, further validating its effectiveness.

2 Related Works

We describe here the existing approaches closely related to our work. For a comprehensive analysis of the state of the art in FL, we refer to [33,44,70].

2.1 Statistical Heterogeneity in Federated Learning

Federated Learning is a topic in continuous growth and evolution. Aiming at a real-world scenario, the non-i.i.d. and unbalanced distribution of users' data poses a significant challenge. The *statistical heterogeneity* of local datasets leads to unstable and slow convergence, suboptimal performance and poor generalization of the global model [73,26,27]. FedAvg [51] defines the standard optimization method and is based on multiple local SGD [56] steps per round. The serverside aggregation is a weighted average of the clients' updates. This simple approach is effective in homogeneous scenarios. Still, it fails to achieve comparable performance against non-i.i.d. data due to local models straying from each other and leading the central model away from the global optimum [34]. To mitigate the effect of the *client drift*, many works enforce regularization in local optimization so that the local model is not led too far apart from the global one [45,34,27,1,43]. Indeed, averaging models/gradients collected from clients having access to a limited subset of tasks may translate into oscillations of the global model and suboptimal performance on the global distribution [48]. Therefore, other lines of research look at improving the aggregation stage using server-side momentum [26] and adaptive optimizers [55], or aggregating task-specific parameters [59,8,9].

In this work, we attempt to explain the behavior of the model in federated scenarios by looking at the loss surface and convergence minima, which is, in our opinion, a fundamental perspective to fully understand the reasons behind the degradation of heterogeneous performance relative to centralized and homogeneous settings. To this end, we focus on explicitly seeking parameters in uniformly low-loss neighborhoods, without any additional communication cost. By encouraging local convergence towards flatter minima, we show that the generalization capacity of the global model is consequently improved. Moreover, thanks to the cyclical average of stochastic weights - accumulated along the trajectory of SGD during rounds on the server-side - broader regions of the weight space are explored, and wider optima are reached. Referring to the terminology introduced by [68], we aim at bridging the participation gap introduced by unseen clients

distributions. Concurrently, [54] provide a theoretical analysis of SAM in FL, matching the convergence rates of the existing methods. Unlike our work, they do not explicitly focus on the issue of statistical heterogeneity in vision tasks.

2.2 Real-world Vision Scenarios in Federated Learning

Research on FL has mainly focused on algorithmic aspects, often overlooking its application to real scenarios and vision tasks. Here, we perform an analysis of the following real-world settings.

Large-scale Classification. Synthetic federated datasets for classification tasks are usually limited in size and do not offer a faithful representation of reality in the data distribution across clients [27]. [27] addresses such issue by adapting the large-scale Google Landmarks v2 [64] to the federated context, using authorship information. We employ the resulting Landmarks-User-160k in our experiments. Semantic Segmentation. A crucial task for real-world applications [19,53], e.g. autonomous driving [58,61], is Semantic Segmentation (SS), which assigns each image pixel to a known category. Most studies of SS in FL focus on medical imaging applications and propose ad hoc techniques to safeguard the patients' privacy [57,46,67,6]. Differently, [52] focuses on object segmentation using prototypical representations. A recently studied application is FL in autonomous driving, motivated by the large amount of privacy-protected data collected by self-driving cars: the authors of [17] propose a new benchmark for analyzing such a scenario, FedDrive. None of those works study the relation between loss landscape and convergence minima of the proposed solution. We apply our approach to the FedDrive benchmark and prove its efficacy in addressing the federated SS task. **Domain Generalization.** When it comes to image data collected from devices around the world, it is realistic to assume there may be different *domains* resulting from the several acquisition devices, light, weather conditions, noise, or viewpoints. With the rising development of FL and the privacy concerns, the problem of Domain Generalization (DG) [7] in a federated setting becomes crucial. DG aims to learn a domain-agnostic model capable of satisfying performances on unseen domains, and its application to federated scenarios is still poorly studied. For instance, [49,62] focus on domain shifts deriving from equipment in the medical field, while [17] analyzes the effects of changing landscapes and weather conditions in the setting of autonomous driving. We show that our approach improves generalization to unseen domains both in classification and SS tasks.

2.3 Flat Minima and Generalization

To understand neural networks' generalization, several theoretical and empirical studies analyze its relationship with the geometry of the loss surface [25,35,16,41,32], connecting sharp minima with poor generalization. "*Flatness*" [25] is defined as the dimension of the region connected around the minimum in which the training loss remains low. Interestingly, it has been shown [32] that sharpnessbased measures highly correlate with generalization performance. The above studies lead to the introduction of Sharpness-Aware Minimization (SAM) [18] which explicitly seeks flatter minima and smoother loss surfaces through a simultaneous minimization of loss sharpness and value during training. As highlighted by [40], SAM is sensitive to parameter re-scaling, weakening the connection between loss sharpness and generalization gap. ASAM [40] solves such issue introducing the concept of adaptive sharpness. Encouraged by their effectiveness across a variety of architectures and tasks[12,4], we ask whether SAM and ASAM can improve generalization in FL as well and find it effective even in the most difficult scenarios. In addition, [20,15] show that local optima found by SGD are connected through a path of near constant loss and that ensambling those points in the weight space leads to high performing networks. Building upon these insights, [29] proposes to average the points traversed by SGD to improve generalization and indeed show the model converges towards wider optima. We modify this approach for FL and use it to cyclically ensemble the models obtained with FedAvg on the server side.

3 Behind the Curtain of Heterogeneous FL

3.1 Federated Learning: Overview

The standard federated framework is based on a central server exchanging messages with K distributed clients. Each device k has access to a privacyprotected dataset \mathcal{D}_k made of N_k images belonging to the input space \mathcal{X} . The goal is to learn a global model f_{θ} parametrized by $\theta \in \mathcal{W} \subseteq \mathbb{R}^d$, where $f_{\theta} : \mathcal{X} \to \mathcal{Y}$ when solving the classification task and $f_{\theta}: \mathcal{X} \to \mathcal{Y}^{N_p}$ in semantic segmentation, with \mathcal{Y} being the output space and N_p the total number of pixels of each image. We assume the structure of θ to be identical across all devices. The learning procedure spans over T communications rounds, during which a subset of clients \mathcal{C} receives the current model parameters θ^t with $t \in [T]$ and trains it on $\mathcal{D}_k \forall k \in \mathcal{C}$, minimizing a local loss function $\mathcal{L}_k(\theta^t) : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}_+$. In FedAvg [51], the global model is updated as a weighted average of the clients' updates θ_k^t , aiming at solving the global objective $\arg \min_{\theta \in \mathbb{R}^d} \frac{1}{N} \sum_{k \in \mathcal{C}} N_k \mathcal{L}_k(\theta)$, with $N = \sum_{k \in \mathcal{C}} N_k$ being the total training images. In particular, from the generalization perspective - defined $\mathcal{D} \triangleq \bigcup_{k \in [K]} \mathcal{D}_k$ the overall clients' data, \mathfrak{D} its distribution and $\mathcal{L}_{\mathcal{D}} =$ $1/\sum_{k} N_k \sum_{k \in [K]} N_k \mathcal{L}_k(\theta)$ the training loss - we aim at learning a model having low population loss $\mathcal{L}_{\mathfrak{D}}(\theta) \triangleq \mathbb{E}_{(x,y)\sim\mathfrak{D}} [\mathbb{E}_{\mathcal{D}}[\mathcal{L}_k(y, f(x, \theta))]]$ [68]. The difference between the population and training losses defines the generalization gap, i.e. the ability of the model to generalize to unseen data [18].

In realistic scenarios, given two clients i and j, \mathcal{D}_i likely follows a different distribution than \mathcal{D}_j , *i.e.* $\mathfrak{D}_i \neq \mathfrak{D}_j$, and the loss $\mathcal{L}_i(\theta) \forall i \in [K]$ is typically non-convex in θ . The loss landscape comprehends a multiplicity of local minima leading to models with different generalization performance, *i.e.* significantly different values of $\mathcal{L}_{\mathfrak{D}}(\theta)$ [18]. Moreover, at each round, the model is likely not to see the entire distribution, further widening the generalization gap [24,23].



Fig. 2: Left: CNN convergence points in distinct federated scenarios with $\alpha \in [0, 0.5, 1k]$ on CIFAR100. Please refer to Appendix C for implementation details. (a) Train loss surface showing the weights obtained at convergence. (b) Test error surface of the same models. Right: Test error surfaces computed on CIFAR100 using three distinct local models after training. (c) When $\alpha = 0$, the local models are not able to generalize to the overall data distribution, being too specialized on the local data. (d) When $\alpha = 1k$, the resulting models are connected through a low-loss region.

3.2 Where Heterogeneous FL Fails at Generalizing

In order to fully understand the behavior of a model trained in a heterogeneous federated scenario, we perform a thorough empirical analysis from different perspectives. Our experimental setup replicates that proposed by [27] both as regards the dataset and the network. The CIFAR100 dataset [39], widely used as benchmark in FL, is split between 100 clients, following a Dirichlet distribution with concentration parameter α . To replicate a heterogeneous scenario, we choose $\alpha \in \{0, 0.5\}$, while α is set to 1000 for the homogeneous one. The model is trained over 20k rounds. Fore more details, please refer to Appendix C.

Model Behavior in Heterogeneous and Homogeneous Scenarios. In Fig. 3, we compare the training trends in centralized, homogeneous and heterogeneous federated settings: in the latter, not only is the trend much noisier and more unstable, but the performance gap is considerable. Consequently, we question the causes of such behavior. First of all, we wonder if the heterogeneous distribution of the data totally inhibits the model from achieving comparable performances: we find it is only a matter of rounds, *i.e.* with a much larger round budget - 10 times larger in our case - the model reaches convergence (Fig. 3). So it becomes obvious the training is somehow slowed down and there is room for improvement. This hypothesis is further validated by the convergence points of the models trained in different settings (Fig. 2): when $\alpha = 1k$ a low-loss region is reached at the end of training, while the same does not happen with lower values of α , meaning that local minima are still to be found. Moreover, the shift between the train and test surfaces suggests us the model trained in the heterogeneous setting $(\alpha = 0)$ is unable to generalize well to unseen data, finding itself in a high-loss region [29]. By analyzing the model behavior, we discover that shifts in client data distribution lead to numerous fluctuations in learning, *i.e.* at each round the model focuses on a subset of the just seen tasks and is unable to generalize to the previously learned ones. This phenomenon is also known as catastrophic interference of neural networks [37] and is typical of the world of multitask learning [10,60]. Fig. 3 highlights this by comparing the accuracy of



Fig. 3: CIFAR100 Accuracy trends. Left: Global model on local distributions with (a) $\alpha = 0$ and (b) 1k @ 20k rounds. Each color represents a local distribution (*i.e.* one class for $\alpha = 0$). (c): $\alpha \in \{0, 0.5, 1k\}$ with necessary rounds to reach convergence.

the global model on the clients' data and the test set when $\alpha = 0$ and $\alpha = 1k$. In the first case, at each round the model achieves very high performances on one class but forgets about the others and this behavior is only slightly attenuated as the training continues. In the homogeneous scenario, on the other hand, the model behaves very similarly on each client and convergence is easily reached, giving way to overfitting as the number of rounds increases.

We analyze the clients' local training for further insights from the characteristics of the updated models. By plotting the position of the weights in the loss landscape after training, we find the models easily overfit the local data distribution (Fig. 2): when tested on the test set, the clients' updates are positioned in very high-error regions and as a result the global model moves away from the minimum, meaning the clients specialize too much on their own data and are not able to generalize to the overall underlying distribution. Moreover, Fig. 2 highlights another relevant issue: models trained on homogeneous distributions are connected through a path of low error and can therefore be ensambled to obtain a more meaningful representation [20], but the same does not hold when $\alpha = 0$, where the models are situated in different loss-value regions. Therefore, FedAvg averages models that are too far apart to lead to a meaningful result.

Federated Training Converges to Sharp Minima. Many works tried to account for this difficulty arising in federated scenarios by enforcing regularization in local optimization not to lead the local model too far apart from the global one [45,34,27,1,43], or by using momentum on the server-side [26], or learning task-specific parameters keeping distinct models on the server-side [59,8,9]. To the best of our knowledge, this is the first work addressing such behavior by looking at the loss landscape. Inspired by a recent trend in Deep Learning connecting the geometry of the loss surface of models trained in non-i.i.d. scenarios with the intention of understanding whether sharp minima may cause the lack of generalization in FL. Following [41], we plot the loss surfaces obtained with models trained in a heterogeneous and in a homogeneous scenario (Fig. 1) showing that both converge to sharp regions, providing a plausible explanation for the highlighted lack of generalization. Additionally, [35] characterizes flatness



Fig. 4: λ_{max}^{k} for each client k as rounds pass

Fig. 5: Hessian eigenspectra of the global model with $\alpha \in \{0, 1k\}$

through the eigenvalues of the Hessian: the dominant eigenvalue λ_{max} evaluates the worst-case landscape curvature, *i.e.* the larger λ_{max} the greater the change in loss in that direction and the steeper the minimum. Hence, we compute the Hessian eigenspectrum (first 50 eigenvalues) using the power iteration mode and analyze it both from the global and local perspectives (Fig. 4.5). Table 1 reports the values of λ_{max} and the ratio λ_{max}/λ_5 , commonly used as a proxy for sharpness [31], as the heterogeneity varies. As expected, λ_{max} is large in all settings when using FedAvg, implying that such method leads the model towards sharp minima regardless of the data distribution, confirming what was noted in the loss landscapes. As for the client-side analysis, we compute the value of λ^k_{max} using the locally updated parameters θ_k^t on the k-th device's data $\mathcal{D}_k \ \forall t \in [T]$. Comparing the i.i.d. and non-i.i.d. settings, we note i) the local values of λ_{max} are much lower if $\alpha = 0$, *i.e.* the clients locally reach wide minima (low Hessian maximum eigenvalue, $\lambda_{max}^k \leq 14$) due to the simplicity of the learned task, *i.e.* a narrow subset of the classes, but the average of the distinct updates drives the model towards sharper minima (high Hessian eigenvalues of the global model, $\lambda_{max} \simeq 94$). ii) When $\alpha \in \{0.5, 1k\}, \lambda_{max}$ decreases as the rounds pass, *i.e.* the global model is moving towards regions with lower curvature, while this is not as evident in the heterogeneous setting. Motivated by these results, we believe that introducing an explicit search for flatter minima can help the model generalize.

4 Seeking Flat Minima in Federated Learning

Common first-order optimizers (*e.g.* SGD [56], Adam [36]) are usually non-robust to unseen data distributions [12], since they only aim at minimizing the

9

Algorithm 1 SAM/ASAM and SWA applied to FedAvg

Rec	quire: Initial random model f^0_{θ} , K clients, T rounds, let size $ \mathcal{B} $ local epoche E cycle length c	arning rates γ_1, γ_2 , neighborhood size $\rho > 0, \eta > 0$, batch
1.	for each round $t = 0$ to $T - 1$ do	
2:	if $t = 0.75 * T$ then	▷ Apply SWA from 75% of training onwards
3:	$\theta_{\text{SWA}} \leftarrow \theta^t$	⊳ Initialize SWA model
4:	end if	
5:	if $t > 0.75 * T$ then	
6:	$\gamma = \gamma(t)$	▷ Compute LR for the round (Eq. 7 in Appendix)
7:	end if	
8:	Subsample a set C of clients	
9:	for each client k in C in parallel do	\triangleright Iterate over subset C of clients
10:	$\theta_{k,0}^{t+1} \leftarrow \theta^t$	
11:	for $e = 0$ to $E - 1$ do	
12:	for $i = 0$ to $N_k/ \mathcal{B} $ do	
13:	Compute gradient $ abla_{m{ heta}} \mathcal{L}_{m{ heta}}(heta_{k,i}^{t+1})$ on batch $m{ heta}$	from \mathcal{D}_k
14:	$\text{Compute } \hat{\epsilon}(\boldsymbol{\theta}_{k,i}^{t+1}) = \rho \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_{k,i}^{t+1}) \big/ \big \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_{k,i}^{t+1}) \big/ \big \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_{k,i}^{t+1}) \big \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_{k,i}^{t+1}) \big $	$ \theta_{k,i}^{t+1}) _2 =: \hat{\epsilon}(\theta) > $ Solve local maximization (Eq. 3)
15:	$\boldsymbol{\theta}_{k,i+1}^{t+1} \leftarrow \boldsymbol{\theta}_{k,i}^{t+1} - \gamma \Big(\left. \left. \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_{k,i}^{t+1}) \right _{\boldsymbol{\theta} + \hat{\boldsymbol{\epsilon}}(\boldsymbol{\theta})} \right. \right.$) \triangleright Local update with sharpness-aware gradient (Eq. 4)
16:	end for	
17:	end for	
18:	Send θ_{1}^{t+1} to the server	
19:	end for	
20:	$ \stackrel{\overline{\theta^{t+1}}}{\leftarrow} \frac{1}{\sum_{k \in \mathcal{C}} N_k} \sum_{k \in \mathcal{C}} N_k \theta_k^{t+1} $	▷ FedAvg
21:	if $t \ge 0.75 * T$ and $mod(t, c) = 0$ then	▷ End of cycle
22:	$n_{\mathrm{models}} \leftarrow t/c$	
23:	$\theta_{\text{SWA}} \leftarrow \frac{\theta_{\text{SWA}} \cdot n_{\text{models}} + \theta^{t+1}}{n_{\text{models}} + 1}$	\triangleright Update SWA average (Eq. 8)
24:	end if	
25:	end for	

training loss $\mathcal{L}_{\mathcal{D}}$, without looking at higher-order information correlating with generalization (e.g. curvature). The federated scenario exacerbates such behavior due to its inherent statistical heterogeneity, resulting in sharp minima and poor generalization. We hypothesize that encouraging the local model to converge towards flatter neighborhoods may help bridging the generalization gap. To this end, we introduce sharpness-aware minimizers, namely SAM [18] and ASAM [40], on the client-side during local training, and Stochastic Weight Averaging [29] on the server-side after the aggregation, adapting the scenario of [29] to FL. By minimizing the sharpness of the loss surface and the generalization gap, the local models are more robust towards unseen data distributions and, when averaged, build a more solid central model. Defined the *sharpness* of a training loss $\mathcal{L}_{\mathcal{D}}$ as $\max_{\|\epsilon\|_p \leq \rho} \mathcal{L}_{\mathcal{D}}(\theta + \epsilon) - \mathcal{L}_{\mathcal{D}}(\theta)$, with ρ being the neighborhood size and $p \in [1, \infty)$, SAM aims at minimizing it by solving $\min_{\theta \in \mathbb{R}^d} \max_{||\epsilon||_p \leq \rho} \mathcal{L}_{\mathcal{D}}(\theta + \epsilon) + \lambda ||\theta||_2^2$. SWA averages weights proposed by SGD, while using a learning rate schedule to explore regions of the weight space corresponding to high performing networks. For a detailed explanation of SAM, ASAM and SWA we refer the reader to Appendix A. Algorithm 1 sums up the details of our approach.

5 Experiments

In this Section, we show the effectiveness of SAM, ASAM and SWA in federated scenarios when addressing tasks of image classification (Sec. 5.1), large-scale

	Algorithm		$\alpha = 0$		α	= 0.5/0.	05	$\alpha = 1000/100$				
	Aigorithin	5cl	10cl	20cl	5cl	10 cl	20 cl	5cl	10 cl	20 cl		
CIEAR100	FedAvg E=1 FedAvg E=2	$30.25 \\ 24.94$	$36.74 \\ 31.81$	$38.59 \\ 35.18$	$40.43 \\ 38.21$	$41.27 \\ 39.59$	$42.17 \\ 40.94$	$49.92 \\ 48.72$	$50.25 \\ 48.64$	$50.66 \\ 48.45$		
	FedSAM FedASAM	31.04 36.04	36.93 39.76	38.56 40.81	44.73 45.61	44.84 46.58	46.05 47.78	54.01 54.81	53.39 54.97	53.97 54.50		
	FedAvg + SWA FedSAM + SWA FedASAM + SWA	39.34 39.30 42.01	39.74 39.51 42.64	39.85 39.24 41.62	43.90 47.96 49.17	44.02 46.76 48.72	42.09 46.47 48.27	50.98 53.90 53.86	50.87 53.67 54.79	50.92 54.36 54.10		
AR10	FedAvg E=1 FedAvg E=2 FedSAM FedASAM	65.00 61.49 70.16 73.66	65.54 62.22 71.09 74.10	68.52 66.36 72.90 76.09	69.24 69.23 73.52 75.61	72.50 69.77 74.81 76.22	73.07 73.48 76.04 76.98	84.46 83.93 84.58 84.77	84.50 84.10 84.67 84.72	84.59 84.21 <u>84.82</u> 84.75		
CIF	FedAvg + SWA FedSAM + SWA FedASAM + SWA	69.71 74.97 76.44	69.54 73.73 75.51	70.19 73.06 76.36	73.48 <u>76.61</u> 76.12	72.80 75.84 76.16	73.81 76.22 76.86	84.35 84.23 84.88	84.32 84.37 84.80	84.47 84.63 84.79		

Table 2: FedSAM, FedASAM and SWA on CIFAR100 and CIFAR10

classification, SS and DG (Sec. 5.2). Their strength indeed lies in finding flatter minima (Sec. 5.1), which consequently help the model to generalize especially in the heterogeneous scenario. We compare our method with algorithms proper of the FL literature and strong data augmentations (Sec. 5.1), commonly used to improve generalization in DL, further validating the efficacy of our proposal. We refer to App. C for implementation details and App. E for the ablation studies.

5.1 The Effectiveness of the Search for Flat Minima in FL

In Sec. 3.2, we have shown that, given a fixed number of rounds, FL models trained in heterogeneous settings present a considerable performance gap compared to their homogeneous counterparts. Indeed, the gap between the two scenarios can be significant with a difference of up to 20% points (Table 2). We identify the clients overspecialization on local data as one of the causes of the poor generalization of the global model to the underlying training distribution. We confirm this by showing the model converges to sharp minima, correlated to a poor generalization capacity. In Table 2, we show that explicitly optimizing for flat minima in both the local training and the server-side aggregation does help improving performances, with evident benefits especially in heterogeneous scenarios. We test SAM, ASAM and their combination with SWA on the federated CIFAR10 and CIFAR100 [39,26,27] with several levels of heterogeneity ($\alpha \in \{0, 0.05, 100\}$ for CIFAR10 and $\alpha \in$ $\{0, 0.5, 1k\}$ for CIFAR100) and clients participation $(K \in \{5, 10, 20\}, i.e. 5\%)$ 10%, 20%). As for CIFAR100, we additionally test our approach on the setting proposed by [55], later referred to as CIFAR100-PAM, where the splits reflect the "coarse" and "fine" label structure proper of the dataset. Since both SAM and ASAM perform a step of gradient ascent and one of gradient descent for each iteration, they should be compared with FedAvg with 2 local epochs. However, the results show FedAvg with E = 2 suffers even more from statistical heterogeneity, so we will compare our baseline with the better-performing FedAvg with E = 1. Our experiments reveal that applying ASAM to FedAvg leads to the best accuracies with a gain of +6% and +8% points respectively on CIFAR100 and CIFAR10 in the most challenging scenario, *i.e.* $\alpha = 0$ and 5 clients per round. This gain

11

·															
		E = 1						E = 2							
Algorithm	Aug	10 clients				20 clients			0 client	s	20 clients				
		@5k	@10k	w/ SWA	@5k	@10k	w/ SWA	@5k	@10k	w/ SWA	@5k	@10k	w/ SWA		
FedAvg		46.60	47.03	52.70	46.51	45.83	50.28	44.58	43.90	51.10	43.31	42.88	47.95		
FedSAM		50.71	53.10	55.44	52.96	53.41	54.67	52.36	52.04	55.23	51.41	51.35	53.41		
FedASAM		49.31	51.10	54.25	47.21	53.50	54.29	49.03	49.33	53.01	53.88	52.94	54.18		
FedAvg	۵.	43.47	49.25	56.71	50.33	49.89	55.74	44.76	46.44	57.15	47.10	47.59	54.40		
FedSAM	ĺnx	42.83	51.92	53.96	49.66	55.77	57.70	42.17	51.04	56.54	53.50	54.75	58.88		
FedASAM	Mi	43.13	51.09	56.31	50.51	52.62	56.89	44.74	50.14	58.31	49.87	50.87	55.86		
FedAvg	4	48.64	48.59	55.40	47.00	46.96	51.70	45.19	45.46	55.40	44.68	44.25	49.39		
FedSAM	nos	48.28	53.53	57.25	52.06	54.37	56.70	49.39	51.88	57.32	52.16	52.37	55.45		
FedASAM	Ę	47.52	52.13	57.01	50.01	50.66	53.54	48.99	50.09	55.77	48.48	48.77	52.00		

Table 3: Accuracy results on CIFAR100-PAM with ResNet18

Table 4: FedAvg, SAM, ASAM and SWA w/ strong data augmentations (Mixup, Cutout)

	Algorithm	SWA	Δ.11 <i>G</i>	$\alpha = 0$			α	= 0.5/0.	05	$\alpha = 1000/100$			
	Algorithm	SWA	Aug	5cl	10cl	20cl	5cl	10cl	20cl	5cl	10cl	20cl	
	FedAvg	×		29.91	33.67	35.67	35.10	37.80	39.34	55.34	55.81	55.98	
	FedSAM	X		30.46	34.10	35.89	38.76	40.31	42.03	54.21	54.94	55.24	
	FedASAM	x	in the	34.04	36.82	36.97	40.71	42.24	44.45	49.75	49.87	49.68	
	FedAvg	~	ix	35.56	36.07	36.08	39.21	39.22	38.31	55.43	55.37	55.39	
8	FedSAM	1	W	35.62	36.25	35.66	42.13	41.95	42.03	52.9	53.14	53.48	
11 11	FedASAM	1		40.08	38.74	37.47	44.53	43.97	44.22	46.97	47.24	46.93	
ΕM	FedAvg	X		24.24	31.55	32.44	37.72	38.45	39.48	53.48	53.83	52.90	
ü	FedSAM	X		23.51	30.92	33.12	40.33	40.31	42.58	54.27	54.75	54.76	
	FedASAM	X	ft	30.05	33.62	34.51	41.86	41.84	43.33	51.88	51.78	53.03	
	FedAvg	1	ft	33.65	34.40	35.03	40.43	40.12	39.32	53.87	54.09	52.75	
	FedSAM	1	ъ	34.00	34.08	34.26	43.09	42.81	42.85	53.78	54.28	53.93	
	FedASAM	1		39.30	37.46	36.27	44.76	43.48	43.95	50.00	49.65	50.81	

is further improved by FedASAM + SWA with a corresponding increase of +12% and +11.5%. The stability introduced by SWA especially helps with lower clients participation, where the trend is noisier. Our ablation studies (Appendix E.3) prove the boost given by SWA is mainly related to the average of the stochastic weights, rather than the cycling learning rate. Table 3 shows the results on CIFAR100-PAM with ResNet18: here SAM and SAM + SWA help more than ASAM.

ASAM and SWA Lead to Flatter Minima in FL. We extend the analysis on the loss landscape and the Hessian eigenspectrum to the models trained with FedSAM, FedASAM and SWA. As expected, both the loss surfaces (Fig. 1) and the Hessian spectra (Fig. 5) indicate us those methods indeed help converging towards flatter minima. The value of λ_{max} goes from 93.5 with FedAvg to 70.3 with FedSAM to 30.1 with FedASAM in the most heterogeneous setting (Table 1). The result is further improved by FedASAM + SWA, obtaining $\lambda_{max} = 24.6$. We notice there is a strict correspondence between the best λ_{max} and the best ratio λ_{max}/λ_5 . Even if the maximum eigenvalue resulting with FedAvg + SWA and FedSAM + SWA is higher than the respective one without SWA, the corresponding lower ratio λ_{max}/λ_5 actually tells us the bulk of the spectrum lies in a lower curvature region [18], proving the effectiveness of SWA. Looking at ASAM's behavior from each client's perspective (Fig. 4), flat minima are achieved from the very beginning of the training and that reflects positively on the model's performance.

ASAM and SWA Enable Strong Data Augmentations in FL. Data augmentations usually play a key role in the performance of a neural network and its ability to generalize [71,65,5], but their design often requires domain expertise

Algorithm	Accu	racy	Absolute	Improvement	Relative Improvement					
Aigoritinii	Centr.	$\alpha = 0$	Centr.	$\alpha = 0$	Centr.	$\alpha = 0$				
SAM	55.22	31.04	+3.02	+0.79	+5.79	+2.61				
ASAM	55.66	36.04	+3.46	+5.79	+6.63	+19.14				
SWA	52.72	39.34	+0.52	+9.09	+1.00	+30.05				
SAM + SWA	55.75	39.30	+0.55	+9.05	+1.06	+29.92				
ASAM + SWA	55.96	42.01	+3.76	+11.76	+7.20	+38.88				
Mixup	58.01	29.91	+5.81	-0.34	+11.13	-1.12				
Cutout	55.30	24.24	+3.10	-6.01	+5.94	-19.87				
Centralized: 52.20 - FedAvg: 30.25										

Table 5: Comparison of improvements (%) in centralized and heterogeneous federated scenarios ($\alpha = 0, 5$ clients) on CIFAR100, computed w.r.t. the reference at the bottom

and greater computational capabilities, two elements not necessarily present in a federated context. In Table 3 and 4, we distinctly apply Mixup [71] and Cutout [14] on CIFAR100-PAM and CIFAR100 (CIFAR10 in Appendix F.2). Surprisingly, both lead to worse performances across all algorithms, so instead of helping the model to generalize, they further slow down training. When combined with our methods, the performance improves in the heterogeneous scenarios w.r.t. the corresponding baseline (FedAvg + data augmentation) and SWA brings a significant boost, enabling the use of data augmentation techniques in FL.

Heterogeneous FL Benefits Even More from Flat Minima. Given the marked improvement brought by SAM, ASAM and their combination with SWA, one might wonder if this simply reflects the gains achieved in the centralized scenario. In Table 5, we prove the positive gap obtained in the heterogeneous federated scenario is larger than the centralized one, showing those approaches are actually helping the training. We also note that while Cutout and Mixup improve the performances in the centralized setting, they do not help in FL, where they achieve a final accuracy worse than FedAvg (Appendix F.1 for $\alpha \in \{0.5, 1k\}$).

Comparison with FL SOTA. We compare our method with FedProx [45], SCAFFOLD [34], FedAvgM [26], FedDyn [1] and AdaBest [63], both on their own and combined with SAM, ASAM and SWA (Table 6). FedProx adds a proximal term to the local objective and, as expected [42,63], does not bring any notable improvement. SCAFFOLD uses control variates to reduce the client drift, exchanging twice the parameters at each round. While performing on par with FedAvg in the homogeneous scenario (84.5% on CIFAR10 and 51.9% on CIFAR100), its performance is heavily affected by the data statistical heterogeneity. The same happens for FedAvgM. FedDyn dynamically aligns global and local stationary points and, as highlighted by [63], is prone to parameters explosion: while it achieves good results on the simpler CIFAR10, it requires heavy gradient clipping and is unable to reach the end of training on CIFAR100. As a solution, AdaBest is proposed, exceeding FedAvg by a few points. Our results demonstrate the consistent effectiveness of FedASAM w.r.t. the SOTA baselines, improving the accuracy by $\approx 6\%$ points on the best SOTA on both datasets. Moreover, by adding ASAM, all FL algorithms notably increase their performance. In particular i) we enable FedAvgM and SCAFFOLD to train in most of the settings with highest heterogeneity, ii) even if limited by the necessary gradient clipping, the results

Algorithm			w/o	SWA			w/ SWA						
		α =	= 0	$\alpha = 0.$	05/0.5	α :	= 0	$\alpha = 0.$	05/0.5				
		5cl	20cl	5cl	20cl	5cl	20cl	5cl	20cl				
CIFAR10	FedAyg FedASAM FedASAM FedAygM FedProx SCAFFOLD FedDyn AdaBest	$\begin{array}{c} 65.00\\ 70.16\\ \textbf{73.66}\\ 10.00\\ 62.72\\ 32.25\\ 67.69\\ 66.77\end{array}$	$\begin{array}{c} 68.52 \\ 72.90 \\ \textbf{76.09} \\ \textbf{10.00} \\ 68.44 \\ 15.56 \\ 73.81 \\ 72.29 \end{array}$	$\begin{array}{r} 69.24\\ 73.52\\ \textbf{75.61}\\ \textbf{10.00}\\ 68.38\\ 54.46\\ 71.36\\ 69.84\end{array}$	73.0776.0476.9878.5173.02 $44.7675.2075.89$	$\begin{array}{r} 69.71 \\ 74.97 \\ 76.44 \\ 10.00 \\ 70.56 \\ 11.98 \\ 77.00 \\ \textbf{78.94} \end{array}$	$\begin{array}{r} 70.19\\ 73.06\\ \textbf{76.36}\\ \textbf{10.00}\\ 70.08\\ \textbf{10.00}\\ \textbf{74.00}\\ 74.00\\ 76.12 \end{array}$	73.48 76.61 76.12 10.00 74.27 33.25 77.99 80.35	73.8176.2276.8684.00 $73.6724.1175.1279.35$				
	FedAvgM + ASAM FedProx + ASAM SCAFFOLD + ASAM FedDyn + SAM AdaBest + ASAM	77.30 73.74 77.78 77.38 77.48	84.89 75.76 77.93 81.00 78.43	77.06 75.32 77.59 79.18 78.41	84.92 77.03 77.80 81.70 79.72	80.88 76.89 75.66 83.81 82.00	85.98 75.92 75.30 86.07 80.80	78.29 76.65 75.32 83.18 81.87	86.03 76.95 75.29 85.57 80.81				
CIFAR100	FedAvg FedASAM FedASAM FedAvgM FedProx SCAFFOLD FedDyn AdaBest FedAvgM + ASAM FedProx + ASAM	$\begin{array}{r} 30.25\\ 31.04\\ \textbf{36.04}\\ 1.00\\ 31.20\\ 1.00\\ 1.00\\ 29.90\\ \hline 1.00\\ 36.10\\ \end{array}$	$\begin{array}{r} 38.59\\ 38.56\\ \textbf{40.81}\\ 40.81\\ 40.64\\ 38.59\\ 1.00\\ 1.40\\ 39.11\\ 39.61\\ 40.91 \end{array}$	$\begin{array}{r} 40.43\\ 44.73\\ 45.61\\ 4.60\\ 39.53\\ 33.26\\ 22.03\\ 36.93\\ \hline 4.60\\ 44.81\end{array}$	$\begin{array}{r} 42.17\\ 46.05\\ 47.78\\ 47.78\\ 42.17\\ 1.00\\ 24.75\\ 43.25\\ \hline 51.65\\ 48.17\\ \end{array}$	$\begin{array}{r} 39.34\\ 39.30\\ 42.01\\ 1.00\\ 39.06\\ 1.00\\ 1.00\\ 44.48\\ \hline 1.00\\ 43.90 \end{array}$	$\begin{array}{r} 39.85\\ 39.24\\ 41.62\\ {\color{red}{53.50}}\\ 39.68\\ 1.00\\ 1.40\\ 44.21\\ {\color{red}{51.58}}\\ 42.06\\ \end{array}$	$\begin{array}{r} 43.90\\ 47.96\\ \textbf{49.17}\\ 4.60\\ 43.98\\ 5.76\\ 8.27\\ 48.20\\ \hline 4.60\\ 48.66\end{array}$	$\begin{array}{r} 42.09\\ 46.47\\ 48.27\\ {\color{red}{53.69}}\\ 41.84\\ {\color{red}{1.00}}\\ 35.15\\ 44.51\\ {\color{red}{56.19}}\\ {\color{red}{48.19}}\end{array}$				
	SCAFFOLD + ASAM FedDyn + ASAM AdaBest + ASAM	$\frac{43.65}{22.16}$	42.61 23.51 45.00	46.50 38.43 45.25	$46.76 \\ 38.60 \\ 49.56$	40.63 17.51 51 75	39.07 19.22	44.87 38.60 51 89	$44.28 \\ 31.06 \\ 51.47$				

Table 6: SOTA comparison on CIFAR10 and CIFAR100 (centralized performance)

reached by FedDyn on CIFAR100 are almost doubled. Lastly, the best results are obtained with ASAM + SWA which stabilizes the noisy learning trends and enables models to converge close to centralized performance with $\alpha = 0$.

5.2 ASAM and SWA in Real World Vision Scenarios

In this Section, we analyze our method in real world scenarios, *i.e.* large scale classification, Semantic Segmentation (SS) for autonomous driving [17] and Domain Generalization (DG) applied to both classification and SS.

Large-scale Classification. We extend our analysis on visual classification tasks to Landmarks-User-160k [27] to validate the effectiveness of SAM, ASAM, and SWA in the presence of real-world challenges such as Non-Identical Class Distribution (different distribution of classes per device), and Imbalanced Client Sizes (varying number of training data per device). Results confirm the benefits of applying client-side sharpness-aware optimizers, especially in combination with server-side weight averaging with an improvement in final accuracy of up to 7%. Semantic Segmentation for Autonomous Driving. SS is a fundamental task for applications of autonomous driving. Due to the private nature of the data collected by self-driving cars, it is reasonable to study this task within a federated scenario. We refer to FedDrive [17] - a new benchmark for autonomous driving in FL - for both settings and baselines. The employed datasets are Cityscapes [13] and IDDA [2] with both uniform and heterogeneous settings. To test the generalization capabilities of the model when facing both semantic and appearance shift, the test domain of IDDA either contains pictures taken in the countryside, or in rainy conditions. The model is tested on both previously seen and unseen domains. As shown in Table 8, ASAM performs best both on Cityscapes and heterogeneous IDDA. The best performance is obtained combining ASAM + SWA with SiloBN [3], keeping the BatchNorm [28] statistics local to each client [47] while sharing the learnable parameters across domains.

Table 7:	Accuracy	Results	(%)	on	Algorithm	Uniform		Cou seen	untry unseen	$_{\rm seen}^{\rm Ra}$	iny unseen		mIoU
Landmark	s-User-160k				FedAvg FedSAM	1		63.31 64.22	48.60 49.74	$\frac{65.16}{64.81}$	$27.38 \\ 30.00$		$ \begin{array}{c} 43.61 \\ 44.58 \end{array} $
	@5k rounds v	v/ SWA 75 w	/ SWA 10	0	FedAvg + SWA			62.74 63.91	$\frac{48.73}{43.28}$	$-64.74 \\ -63.24$	$\frac{31.32}{47.72}$	1	45.64
FedAvg	61.91	66.05	67.52		FedSAM + SWA FedASAM + SWA	1	~	$\begin{array}{c} 62.26 \\ 60.78 \end{array}$	46.26 44.23	63.69 63.18	48.40 51.76	VPES	45.29 45.69
FedASAM	$63.72 \\ 64.23$	$67.11 \\ 67.17$	68.32		FedAvg FedSAM	×	q	$\frac{42.06}{43.28}$	36.04 37.83	$39.50 \\ 39.65$	$\frac{24.59}{29.27}$	'YSC/	$\frac{38.65}{41.22}$
Centralized	1	74.03			FedASAM FedAvg + SWA	x		43.67 37.16	$\frac{36.11}{37.48}$	$\frac{41.68}{37.06}$	30.07 42.33	5	42.27 42.48
					FedSAM + SWA	X		44.26	$\frac{40.45}{20.72}$	38.15	45.25		43.42
					SiToBN - SWA	x		45.86	$\frac{39.12}{32.77}$	48.09	39.67	1	45.96^{-1}
					SiloBN + SAM SiloBN + ASAM	×		$\frac{46.88}{46.57}$	33.71 35.22	48.22 48.33	40.08 40.76		49.10 49.75

Table 8: Federated SS on Cityscapes and IDDA. Results in mIoU (%) @ 1.5k rounds

Domain Generalization. To further show the generalization performance acquired by the model trained with SAM, ASAM and SWA, we test it on the corrupted CIFAR datasets [24]. The test images are altered by 19 corruptions each with 5 levels of severity. Fig. 6 shows the results on the highest severity and once again validate the efficacy of seeking flat minima in FL (complete results in App. D).



Fig. 6: Domain generalization in FL. Results with $\alpha = 0$, 20 clients, severity level 5.

6 Conclusions

Heterogeneous Federated Learning suffers from degraded performances and slowdown in training due to the poor generalization of the learned global model. Inspired by recent trends in deep learning connecting the loss landscape and the generalization gap, we analyzed the behavior of the model through the lens of the geometry of the loss surface and linked the lack of generalization to convergence towards sharp minima. As a solution, we introduced Sharpness-Aware Minimization, its adaptive version and Stochastic Weight Averaging in FL for encouraging convergence towards flatter minima. We showed the effectiveness of this approach in several vision tasks and datasets.

Acknowledgments. We thank L. Fantauzzo for her help with the SS experiments. We acknowledge the CINECA HPC infrastructure. Work funded by CINI.

References

- Acar, D.A.E., Zhao, Y., Navarro, R.M., Mattina, M., Whatmough, P.N., Saligrama, V.: Federated learning based on dynamic regularization. International Conference on Learning Representations (2021) 2, 3, 7, 12
- Alberti, E., Tavera, A., Masone, C., Caputo, B.: Idda: a large-scale multi-domain dataset for autonomous driving. IEEE Robotics and Automation Letters 5(4), 5526–5533 (2020) 13
- Andreux, M., Terrail, J.O.d., Beguier, C., Tramel, E.W.: Siloed federated learning for multi-centric histopathology datasets. In: Domain Adaptation and Representation Transfer, and Distributed and Collaborative Learning, pp. 129–139. Springer (2020) 13
- Bahri, D., Mobahi, H., Tay, Y.: Sharpness-aware minimization improves language model generalization. arXiv preprint arXiv:2110.08529 (2021) 5
- Bello, I., Fedus, W., Du, X., Cubuk, E.D., Srinivas, A., Lin, T.Y., Shlens, J., Zoph, B.: Revisiting resnets: Improved training and scaling strategies. Advances in Neural Information Processing Systems 34 (2021) 11
- Bercea, C.I., Wiestler, B., Rueckert, D., Albarqouni, S.: Feddis: Disentangled federated learning for unsupervised brain pathology segmentation. arXiv preprint arXiv:2103.03705 (2021) 4
- Blanchard, G., Lee, G., Scott, C.: Generalizing from several related classification tasks to a new unlabeled sample. Advances in neural information processing systems 24 (2011) 3, 4
- Briggs, C., Fan, Z., Andras, P.: Federated learning with hierarchical clustering of local updates to improve training on non-iid data. In: 2020 International Joint Conference on Neural Networks (IJCNN). pp. 1–9. IEEE (2020) 3, 7
- Caldarola, D., Mancini, M., Galasso, F., Ciccone, M., Rodolà, E., Caputo, B.: Cluster-driven graph federated learning over multiple domains. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop. pp. 2749–2758 (2021) 3, 7
- 10. Caruana, R.: Multitask learning. Machine learning 28(1), 41-75 (1997) 6
- Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence 40(4), 834–848 (2017) 3
- Chen, X., Hsieh, C.J., Gong, B.: When vision transformers outperform resnets without pre-training or strong data augmentations. In: International Conference on Learning Representations (2022) 2, 5, 8
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3213–3223 (2016) 13
- 14. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017) 12
- Draxler, F., Veschgini, K., Salmhofer, M., Hamprecht, F.: Essentially no barriers in neural network energy landscape. In: International conference on machine learning. pp. 1309–1318. PMLR (2018) 5
- Dziugaite, G.K., Roy, D.M.: Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. arXiv preprint arXiv:1703.11008 (2017) 2, 4, 7

- 16 D. Caldarola et al.
- Fantauzzo, L., Fani', E., Caldarola, D., Tavera, A., Cermelli, F., Ciccone, M., Caputo, B.: Feddrive: Generalizing federated learning to semantic segmentation in autonomous driving. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (2022) 4, 13
- Foret, P., Kleiner, A., Mobahi, H., Neyshabur, B.: Sharpness-aware minimization for efficiently improving generalization. International Conference on Learning Representations (2021) 2, 4, 5, 9, 11
- Garcia-Garcia, A., Orts-Escolano, S., Oprea, S., Villena-Martinez, V., Garcia-Rodriguez, J.: A review on deep learning techniques applied to semantic segmentation. arXiv preprint arXiv:1704.06857 (2017) 4
- Garipov, T., Izmailov, P., Podoprikhin, D., Vetrov, D.P., Wilson, A.G.: Loss surfaces, mode connectivity, and fast ensembling of dnns. Advances in neural information processing systems **31** (2018) 5, 7
- Gong, X., Sharma, A., Karanam, S., Wu, Z., Chen, T., Doermann, D., Innanje, A.: Ensemble attention distillation for privacy-preserving federated learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 15076–15086 (October 2021) 1
- Guo, P., Wang, P., Zhou, J., Jiang, S., Patel, V.M.: Multi-institutional collaborations for improving deep learning-based magnetic resonance image reconstruction using federated learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2423–2432 (June 2021) 1
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., et al.: The many faces of robustness: A critical analysis of out-of-distribution generalization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8340–8349 (2021) 5
- Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. International Conference on Learning Representations (2019) 5, 14
- Hochreiter, S., Schmidhuber, J.: Flat minima. Neural computation 9(1), 1–42 (1997)
 2, 4
- Hsu, T.M.H., Qi, H., Brown, M.: Measuring the effects of non-identical data distribution for federated visual classification. NeurIPS Workshop (2019) 1, 3, 7, 10, 12
- Hsu, T.M.H., Qi, H., Brown, M.: Federated visual classification with real-world data distribution. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16. pp. 76–92. Springer (2020) 1, 3, 4, 6, 7, 10, 13
- Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: International conference on machine learning. pp. 448–456. PMLR (2015) 13
- Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D., Wilson, A.G.: Averaging weights leads to wider optima and better generalization. Uncertainty in Artificial Intelligence (UAI) (2018) 2, 5, 6, 7, 9
- Jastrzebski, S., Kenton, Z., Ballas, N., Fischer, A., Bengio, Y., Storkey, A.: On the relation between the sharpest directions of dnn loss and the sgd step length. International Conference on Learning Representations (2019) 2
- Jastrzebski, S., Szymczak, M., Fort, S., Arpit, D., Tabor, J., Cho, K., Geras, K.: The break-even point on optimization trajectories of deep neural networks. arXiv preprint arXiv:2002.09572 (2020) 8

Improving Generalization in Federated Learning by Seeking Flat Minima

- Jiang, Y., Neyshabur, B., Mobahi, H., Krishnan, D., Bengio, S.: Fantastic generalization measures and where to find them. arXiv preprint arXiv:1912.02178 (2019) 2, 4, 7
- Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al.: Advances and open problems in federated learning. arXiv preprint arXiv:1912.04977 (2019) 3
- Karimireddy, S.P., Kale, S., Mohri, M., Reddi, S., Stich, S., Suresh, A.T.: Scaffold: Stochastic controlled averaging for federated learning. In: International Conference on Machine Learning. pp. 5132–5143. PMLR (2020) 2, 3, 7, 12
- Keskar, N.S., Mudigere, D., Nocedal, J., Smelyanskiy, M., Tang, P.T.P.: On largebatch training for deep learning: Generalization gap and sharp minima. International Conference on Learning Representations (2017) 2, 4, 7
- 36. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. ICLR (2015) 8
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al.: Overcoming catastrophic forgetting in neural networks. Proceedings of the national academy of sciences 114(13), 3521–3526 (2017) 6
- Kleinberg, B., Li, Y., Yuan, Y.: An alternative view: When does sgd escape local minima? In: International Conference on Machine Learning. pp. 2698–2707. PMLR (2018) 2
- Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009) 6, 10
- Kwon, J., Kim, J., Park, H., Choi, I.K.: Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. International Conference on Machine Learning (2021) 2, 5, 7, 9
- 41. Li, H., Xu, Z., Taylor, G., Studer, C., Goldstein, T.: Visualizing the loss landscape of neural nets. In: Neural Information Processing Systems (2018) 2, 4, 7
- Li, Q., Diao, Y., Chen, Q., He, B.: Federated learning on non-iid data silos: An experimental study. arXiv preprint arXiv:2102.02079 (2021) 12
- Li, Q., He, B., Song, D.: Model-contrastive federated learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10713–10722 (2021) 1, 3, 7
- Li, T., Sahu, A.K., Talwalkar, A., Smith, V.: Federated learning: Challenges, methods, and future directions. IEEE Signal Processing Magazine 37(3), 50–60 (2020) 3
- Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A., Smith, V.: Federated optimization in heterogeneous networks. Proceedings of Machine Learning and Systems 2, 429–450 (2020) 1, 2, 3, 7, 12
- 46. Li, W., Milletarì, F., Xu, D., Rieke, N., Hancox, J., Zhu, W., Baust, M., Cheng, Y., Ourselin, S., Cardoso, M.J., et al.: Privacy-preserving federated brain tumour segmentation. In: International workshop on machine learning in medical imaging. pp. 133–141. Springer (2019) 4
- Li, Y., Wang, N., Shi, J., Liu, J., Hou, X.: Revisiting batch normalization for practical domain adaptation. ICLR Workshop (2017) 13
- Lin, T., Kong, L., Stich, S.U., Jaggi, M.: Ensemble distillation for robust model fusion in federated learning. arXiv preprint arXiv:2006.07242 (2020) 3
- Liu, Q., Chen, C., Qin, J., Dou, Q., Heng, P.A.: Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1013–1023 (2021) 1, 4

- 18 D. Caldarola et al.
- Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3431–3440 (2015) 3
- McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communicationefficient learning of deep networks from decentralized data. In: Artificial intelligence and statistics. pp. 1273–1282. PMLR (2017) 1, 3, 5
- Michieli, U., Ozay, M.: Prototype guided federated learning of visual feature representations. arXiv preprint arXiv:2105.08982 (2021) 4
- Ouahabi, A., Taleb-Ahmed, A.: Deep learning for real-time semantic segmentation: Application in ultrasound imaging. Pattern Recognition Letters 144, 27–34 (2021)
 4
- Qu, Z., Li, X., Duan, R., Liu, Y., Tang, B., Lu, Z.: Generalized federated learning via sharpness aware minimization. In: Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., Sabato, S. (eds.) Proceedings of the 39th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 162, pp. 18250–18280. PMLR (17–23 Jul 2022) 4
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., McMahan, H.B.: Adaptive federated optimization. International Conference on Learning Representations (2021) 3, 10
- 56. Ruder, S.: An overview of gradient descent optimization algorithms. arXiv preprint arXiv:1609.04747 (2016) 3, 8
- 57. Sheller, M.J., Reina, G.A., Edwards, B., Martin, J., Bakas, S.: Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation. In: International MICCAI Brainlesion Workshop. pp. 92–104. Springer (2018) 4
- Siam, M., Gamal, M., Abdel-Razek, M., Yogamani, S., Jagersand, M., Zhang, H.: A comparative study of real-time semantic segmentation for autonomous driving. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 587–597 (2018) 4
- Smith, S.L., Le, Q.V.: A bayesian perspective on generalization and stochastic gradient descent. International Conference on Learning Representations (2018) 2, 3, 7
- Smith, V., Chiang, C.K., Sanjabi, M., Talwalkar, A.S.: Federated multi-task learning. Advances in neural information processing systems **30** (2017) 6
- 61. Tavera, A., Cermelli, F., Masone, C., Caputo, B.: Pixel-by-pixel cross-domain alignment for few-shot semantic segmentation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1626–1635 (2022) 4
- Tian, C.X., Li, H., Wang, Y., Wang, S.: Privacy-preserving constrained domain generalization for medical image classification. arXiv preprint arXiv:2105.08511 (2021) 4
- Varno, F., Saghayi, M., Rafiee, L., Gupta, S., Matwin, S., Havaei, M.: Minimizing client drift in federated learning via adaptive bias estimation. arXiv preprint arXiv:2204.13170 (2022) 12
- 64. Weyand, T., Araujo, A., Cao, B., Sim, J.: Google landmarks dataset v2-a largescale benchmark for instance-level recognition and retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2575–2584 (2020) 4
- Xie, C., Tan, M., Gong, B., Wang, J., Yuille, A.L., Le, Q.V.: Adversarial examples improve image recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 819–828 (2020) 11

Improving Generalization in Federated Learning by Seeking Flat Minima

- Yao, C.H., Gong, B., Qi, H., Cui, Y., Zhu, Y., Yang, M.H.: Federated multi-target domain adaptation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1424–1433 (2022) 1
- Yi, L., Zhang, J., Zhang, R., Shi, J., Wang, G., Liu, X.: Su-net: an efficient encoderdecoder model of federated learning for brain tumor segmentation. In: International Conference on Artificial Neural Networks. pp. 761–773. Springer (2020) 4
- Yuan, H., Morningstar, W., Ning, L., Singhal, K.: What do we mean by generalization in federated learning? NeurIPS Workshop (2021) 3, 5
- Yue, X., Nouiehed, M., Kontar, R.A.: Salr: Sharpness-aware learning rates for improved generalization. arXiv preprint arXiv:2011.05348 (2020) 2
- Zhang, C., Xie, Y., Bai, H., Yu, B., Li, W., Gao, Y.: A survey on federated learning. Knowledge-Based Systems 216, 106775 (2021) 3
- Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. International Conference on Learning Representations (2018) 11, 12
- Zhang, L., Luo, Y., Bai, Y., Du, B., Duan, L.Y.: Federated learning for non-iid data via unified feature learning and optimization objective alignment. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 4420– 4428 (October 2021) 1
- Zhao, Y., Li, M., Lai, L., Suda, N., Civin, D., Chandra, V.: Federated learning with non-iid data. arXiv preprint arXiv:1806.00582 (2018) 1, 3