

Monge SAM: Robust Reparameterization-Invariant Sharpness-Aware Minimization Based on Loss Geometry

Albert Kjøller Jacobsen¹ Georgios Arvanitidis¹

Abstract

Recent studies on deep neural networks show that flat minima of the loss landscape correlate with improved generalization. Sharpness-aware minimization (SAM) efficiently finds flat regions by updating the parameters according to the gradient at an adversarial perturbation. The perturbation depends on the Euclidean metric, making SAM non-invariant under reparameterizations, which blurs sharpness and generalization. We propose Monge SAM (M-SAM), a reparameterization invariant version of SAM by considering a Riemannian metric in the parameter space induced naturally by the loss surface. Compared to previous approaches, M-SAM works under any modeling choice, relies only on mild assumptions while being as computationally efficient as SAM. We theoretically argue that M-SAM varies between SAM and gradient descent (GD), which increases robustness to hyperparameter selection and reduces attraction to suboptimal equilibria like saddle points. We demonstrate this behavior both theoretically and empirically on a multi-modal representation alignment task.

1. Introduction

What makes overparameterized deep neural networks capable of generalizing as well as they do? Inspired by the early work of Hochreiter & Schmidhuber (1997) that argued how flat regions of the parameter space correspond to networks with low expected overfitting, recent works showed correlation between flatness and generalization capabilities (Keskar et al., 2016; Dziugaite & Roy, 2017; Izmailov et al., 2018; Jastrzebski et al., 2018; Chaudhari et al., 2019; Jiang et al., 2019), including downstream performance of pre-trained large language models (LLMs) (Liu et al., 2023).

¹Section for Cognitive Systems, DTU Compute, Technical University of Denmark, Kongens Lyngby, Denmark. Correspondence to: Albert Kjøller Jacobsen <akjja@dtu.dk>.

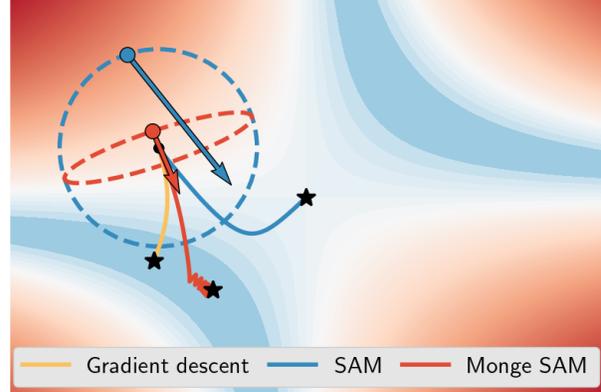


Figure 1. The SAM finds the adversarial perturbation within a Euclidean ball (---) which upper bounds the M-SAM perturbation that is based on the local geometry of the loss (---), implying an adaptive trade-off between SAM and GD. In a loss defined by $\ell(\theta) = (1 - \theta_1\theta_2)^2$ with banana-shaped minima at $\theta_1 = 1/\theta_2$, M-SAM is less prone to get attracted to the saddle point at $\theta_s = (0, 0)$ than SAM. M-SAM can reach lower losses like GD while being capable of walking along minima, eventually finding the flattest global minimum at $\theta_{\text{flat}}^* = (-1, -1)$. We run 200 steps from $\theta_0 = (-\frac{3}{2}, \frac{1}{2})$ with $\rho = 1$ and a learning rate of 0.01. Arrows represent the respective gradients (rescaled) at the perturbed points.

The sharpness-aware minimization (SAM) approach to optimization (Foret et al., 2020) has increasingly gained traction due to its practical formulation that seeks flat parameter space regions without requiring information about the curvature through the Hessian, as it only relies on first-order derivatives. Though SAM might enhance generalization, it is known that sharp minima can also generalize equivalently whenever the learned function is reparameterized accordingly (Dinh et al., 2017; Kristiadi et al., 2024). Therefore, several extensions to SAM have been proposed; Kwon et al. (2021) proposed adaptive SAM (ASAM) to partly tackle this reparameterization issue by making SAM scale-invariant, where Fisher SAM (Kim et al., 2022) adapts the perturbation step of SAM by employing a Riemannian metric induced in the parameter space under the Information Geometry framework, thereby respecting the local geometry of probabilistic models. Similarly, Riemannian SAM (Yun

& Yang, 2024) generalizes this framework, which includes Fisher SAM, by considering the parameter space of a deep neural network as a predefined Riemannian manifold, as a sphere, and computing perturbations and gradient updates on the manifold.

In this paper, we propose a novel approach to tackle the reparameterization issue of SAM. We naturally induce a Riemannian metric in the parameter space that captures the geometry of the loss surface, namely the Monge metric, and thereby constrain the region in which the adversarial SAM perturbation is searched for. Our approach is more general than previously proposed approaches as it does not rely on a probabilistic formulation or predefined manifolds while still respecting the local structure of the loss. Due to the simplicity of this metric, our proposed method has an analytical expression of the adversarial perturbation, in contrast to Fisher SAM, which typically requires regularizing and approximating the metric through diagonalization.

Our main contributions include the following:

1. We establish *Monge SAM (M-SAM)*; a novel sharpness-aware minimizer that exploits the Monge metric.
2. M-SAM relies on mild assumptions and is not based on probabilistic model formulations or predefining the parameter manifold; it works for any modeling choice.
3. We theoretically justify that M-SAM is less prone to get attracted to suboptimal equilibria like saddle points. We extend on previous studies by analyzing stability of the SAM gradient flow using perturbation theory.
4. We provide empirical evidence of M-SAM being more robust to hyperparameter selection than SAM, having a beneficial impact on performance for some tasks.

2. Background

Notation. We denote the K -dimensional parameter space of a parametric data-dependent model, $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^C$ (e.g. a neural network) by $\theta = (\theta_1, \dots, \theta_K) \in \Theta \subseteq \mathbb{R}^K$. We define general loss functions as $\ell : \Theta \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ where $\mathcal{X} \in \mathbb{R}^D$ and $\mathcal{Y} \in \mathbb{R}^C$ are the input and output spaces, respectively. For simplicity we disregard the dependency on the data $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and denote the general loss function by $\ell(\theta)$ for a specific realization of the model f_{θ} .

2.1. Sharpness-Aware Minimization

As mentioned, SAM (Foret et al., 2020) seeks flat minima without relying on Hessian-based measures as such would be impractical to compute. SAM seeks to minimize the loss at a perturbed parameter set, $\theta + \delta$, where δ is the adversarial

perturbation. Thus SAM’s optimization objective is:

$$\min_{\theta} \max_{\|\delta\|_{\mathbf{M}} \leq \rho} \ell(\theta + \delta), \quad (1)$$

where $\|\delta\|_{\mathbf{M}}^2 = \delta^{\top} \mathbf{M} \delta$ defines the local norm under the positive-definite matrix \mathbf{M} .

The type of flatness minimized by SAM appears in the objective by adding and subtracting the regular loss term and rearranging gives:

$$\min_{\theta} \underbrace{\ell(\theta)}_{\text{Loss term}} + \underbrace{\max_{\|\delta\|_{\mathbf{M}} \leq \rho} \ell(\theta + \delta) - \ell(\theta)}_{\text{Sharpness term}} \quad (2)$$

Thus, SAM’s notion of sharpness is the difference in loss values evaluated at the original parameters and the worst-case perturbation of the parameters, under the local norm constraint. In practice, finding the worst-case perturbation is handled by solving the dual-norm problem that occurs when approximating the objective with a first-order Taylor expansion. Previous works consider $\mathbf{M} = \mathbb{I}_K$ as the default choice, i.e. searching within the Euclidean ball, for which the worst-case perturbation is the rescaled gradient:

$$\delta_{\text{SAM}}^* = \frac{\rho}{\|\nabla \ell(\theta)\|_2} \cdot \nabla \ell(\theta). \quad (3)$$

After computing the worst-case perturbation, minimizing the SAM objective requires an additional backward pass under the *base optimizer* (e.g. stochastic gradient descent). Remark that the scaling factor $\tilde{\rho}_{\text{SAM}} := \rho / \|\nabla \ell(\theta)\|_2$ will later be referred to as the *effective perturbation size*.

2.2. Fisher SAM

Fisher SAM (Kim et al., 2022) exploits the fact that the KL-divergence between two infinitesimally close parametric distributions can be approximated by the norm under a Riemannian metric. In particular, an approximation to this metric is the empirical Fisher information matrix. Leveraging this fact, the optimization objective is rephrased as

$$\min_{\theta} \max_{\delta^{\top} \mathbf{F}(\theta) \delta \leq \rho^2} \ell(\theta + \delta). \quad (4)$$

Approximating this with a first-order Taylor expansion leads to a quadratically constrained linear program for the worst-case perturbation. Solving the Lagrangian gives:

$$\delta_{\text{Fisher}}^* = \rho \cdot \frac{\mathbf{F}(\theta)^{-1} \nabla \ell(\theta)}{\sqrt{\nabla \ell(\theta)^{\top} \mathbf{F}(\theta)^{-1} \nabla \ell(\theta)}}. \quad (5)$$

Fisher SAM is reparameterization-invariant and respects the local geometry of the probabilistic model through the Fisher metric. The Fisher metric is defined as the sum of the outer products of the gradient of log-likelihoods,

$\mathbf{F}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\theta} [\nabla \log p(y|\mathbf{x}, \boldsymbol{\theta}) \nabla \log p(y|\mathbf{x}, \boldsymbol{\theta})^\top]$, which grows quadratically with the number of parameters, and reveals that Fisher SAM only works for probabilistic models. In practice, the implementation relies on approximation by mini-batching and diagonalization as $\hat{\mathbf{F}}(\boldsymbol{\theta}) = \text{diag} \left[\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla \log p(y_i|\mathbf{x}_i, \boldsymbol{\theta}) \right]^2$ with \mathcal{S} being a batch. This ad-hoc approximation is necessary as inverting the original metric is prohibitively expensive, while to avoid division by zero the inverse is defined as $\hat{f}_i^{-1} = 1/(1 + \eta f_i)$ where f_i denotes the diagonal elements of $\hat{\mathbf{F}}(\boldsymbol{\theta})$ and η a hyperparameter. Note that diagonalization removes correlation, which underrepresents the local structure.

3. Monge SAM

We develop a geometry-aware version of SAM that is not based on probabilistic model formulations or predefined parameter manifolds. Intuitively, we approach this problem by modifying the parameter space Θ to locally be aware of the training loss structure. We propose equipping the parameter space with a simple Riemannian metric that encodes the loss surface geometry, which is computationally efficient and enables us to obtain a closed-form expression of the worst-case perturbation.

3.1. The Monge metric

Given the definition of a general loss function $\ell(\boldsymbol{\theta})$, we consider the loss surface of the model $f_{\boldsymbol{\theta}}$ as a K -dimensional smooth manifold embedded in \mathbb{R}^{K+1} as follows $\mathcal{M} = g(\boldsymbol{\theta}) = [\theta_1, \dots, \theta_K, \ell(\boldsymbol{\theta})] \in \mathbb{R}^{d+1}$. The parameter space Θ represents the *intrinsic coordinates* of the manifold \mathcal{M} , and at a point $g(\boldsymbol{\theta}) \in \mathcal{M}$ the tangent space $\mathcal{T}_{g(\boldsymbol{\theta})}\mathcal{M}$ is spanned by the Jacobian $\mathbf{J}_g : \Theta \rightarrow \mathbb{R}^{K+1 \times K}$. We can therefore write a tangent vector to the manifold in terms of its intrinsic coordinates $\vec{v} \in \mathbb{R}^K$ as $\mathbf{J}_g(\boldsymbol{\theta}) \vec{v}$. Likewise, the inner product between two tangential vectors, \vec{v}_1 and \vec{v}_2 , in the same tangent plane is $\langle \mathbf{J}_g(\boldsymbol{\theta}) \vec{v}_1, \mathbf{J}_g(\boldsymbol{\theta}) \vec{v}_2 \rangle = \vec{v}_1^\top \mathbf{J}_g(\boldsymbol{\theta})^\top \mathbf{J}_g(\boldsymbol{\theta}) \vec{v}_2$. Evidently, the induced *Riemannian metric* $\mathbf{G}(\boldsymbol{\theta}) = \mathbf{J}_g(\boldsymbol{\theta})^\top \mathbf{J}_g(\boldsymbol{\theta})$ provides a notion of local inner products on the manifold through its intrinsic coordinates and captures the local geometry.

Due to the parametrization $g(\boldsymbol{\theta})$ of the manifold \mathcal{M} the Jacobian takes a simple form, i.e. $\mathbf{J}_g(\boldsymbol{\theta}) = [\mathbb{I}_K, \nabla \ell(\boldsymbol{\theta})]^\top$. Consequently, this gives a simple expression for the metric:

$$\mathbf{G}(\boldsymbol{\theta}) = \mathbb{I}_K + \nabla \ell(\boldsymbol{\theta}) \nabla \ell(\boldsymbol{\theta})^\top, \quad (6)$$

which is known as the *Monge* metric and consists of the outer product of the loss gradient with itself, regularized by adding 1's to the diagonal entries. The inherent regularization ensures positive definiteness of $\mathbf{G}(\boldsymbol{\theta})$ and thereby invertibility. Following the Sherman-Morrison formula the

inverse metric takes the form:

$$\mathbf{G}(\boldsymbol{\theta})^{-1} = \mathbb{I}_K - \frac{\nabla \ell(\boldsymbol{\theta}) \nabla \ell(\boldsymbol{\theta})^\top}{1 + \|\nabla \ell(\boldsymbol{\theta})\|_2^2}. \quad (7)$$

Using the Monge metric rather than the Fisher metric has several advantages; while gradients of the log-likelihood can indeed be obtained for probabilistic models, accessing these quantities requires careful implementation when relying on automatic differentiation frameworks. Contrarily, the gradient of the loss is straight-forward to obtain for any implementation. Secondly, the inherent regularization ensures a closed-form expression of the inverse Monge metric, thereby allowing us to use its full expressivity without using tricks like diagonalization and explicit regularization.

3.2. The Monge SAM perturbation

Leveraging the notion of local inner products as expressed using the Monge metric, we slightly change the objective:

$$\min_{\boldsymbol{\theta}} \max_{\delta^\top \mathbf{G}(\boldsymbol{\theta}) \delta \leq \rho^2} \ell(\boldsymbol{\theta} + \delta). \quad (8)$$

Though solving for the worst-case perturbation gives a similar expression as in Equation 5, this further simplifies due to the inverse Monge metric being analytically accessible, resulting in a closed-form expression for M-SAM's worst-case perturbation:

$$\delta_{\text{M-SAM}}^* = \frac{1}{\sqrt{1 + \|\nabla \ell(\boldsymbol{\theta})\|_2^2}} \cdot \underbrace{\frac{\rho}{\|\nabla \ell(\boldsymbol{\theta})\|_2} \cdot \nabla \ell(\boldsymbol{\theta})}_{= \tilde{\delta}_{\text{SAM}}^*}. \quad (9)$$

We see that the geometry-aware worst-case perturbation is a rescaled version of SAM's worst-case perturbation and define $\tilde{\rho}_{\text{M-SAM}} := \rho / (\|\nabla \ell(\boldsymbol{\theta})\|_2 \cdot \sqrt{1 + \|\nabla \ell(\boldsymbol{\theta})\|_2^2})$ as M-SAM's effective perturbation radius. As such, the computational requirements are comparable to SAM.

Monge SAM is conservative. We now consider how the adjusted worst-case perturbation behaves under two extreme cases, namely at locations in the parameter space where the loss is either 1) very steep or 2) close to being a stationary point, i.e. $\|\nabla \ell(\boldsymbol{\theta})\|_2 \rightarrow +\infty$ and $\|\nabla \ell(\boldsymbol{\theta})\|_2 \rightarrow 0$, respectively. The asymptotic behavior is:

$$\begin{aligned} \lim_{\|\nabla \ell(\boldsymbol{\theta})\|_2 \rightarrow +\infty} \delta_{\text{M-SAM}}^* &= \mathbf{0} \\ \lim_{\|\nabla \ell(\boldsymbol{\theta})\|_2 \rightarrow 0} \delta_{\text{M-SAM}}^* &= \tilde{\delta}_{\text{SAM}}^*. \end{aligned}$$

Put differently, the effective perturbation radius of M-SAM is upper bounded by that of SAM, i.e. $\tilde{\rho}_{\text{M-SAM}} \leq \tilde{\rho}_{\text{SAM}}$. Remarkably, M-SAM appears to *trade-off* between gradient descent (GD) and SAM behavior by restricting the perturbation inversely to the growth of the training loss where the original SAM perturbation can in principle explode.

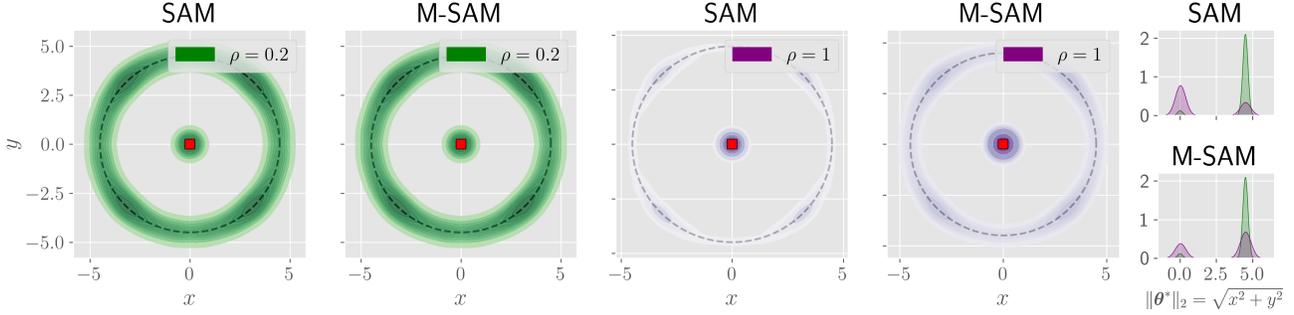


Figure 2. Attraction to maxima. We consider the scaled 2D sinc-function given by $\text{sinc}(x, y) = 5 \cdot \sin(x^2 + y^2) / (x^2 + y^2)$ and draw $N = 200$ samples of $\theta = (x, y)$, uniformly distributed within the centered unit square. For each sample, we run SAM and M-SAM for 200 steps with a learning rate of 0.05 and $\rho \in \{0.2, 1\}$ and plot the distribution of the converged estimates, θ^* . The larger perturbation radius $\rho = 1$ makes the global maximum (■) at $\theta = (0, 0)$ a stronger attractor for SAM, as for $\rho = 0.2$ SAM is more likely to descend into the surrounding circular minima range (---). We observe similar trends for M-SAM, yet see that the conservative property restricts how strong attractor the maxima is, even for the high perturbation radius of $\rho = 1$.

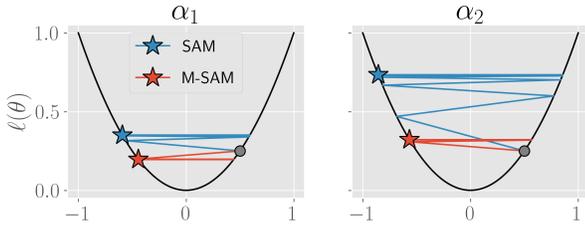


Figure 3. Conservative property. SAM vs. M-SAM behavior in a simple paraboloid loss given by $\ell(\theta) = \theta^2$ for two learning rates, $\alpha_1 < \alpha_2$ when fixing the perturbation radius $\rho = 0.3$ and taking 11 steps. M-SAM always converges to lower values than SAM and is additionally limitedly affected by large learning rates, revealing M-SAM’s conservativeness that originates from being loss-aware.

In Figure 3 we show how M-SAM is biased towards GD behavior. We apply SAM and M-SAM in a symmetric 1D toy loss defined by a second-order polynomial, i.e. $\ell(\theta) = \theta^2$. We compare trajectories of the two methods for two learning rates $\alpha_1 = 0.66$ and $\alpha_2 = 0.74$ when fixing the perturbation radius to $\rho = 0.3$. For α_1 the two methods behave differently; the SAM perturbation - and consequently also the gradient update - is sufficiently large for SAM to diverge from the initial location whereas M-SAM has a lower effective perturbation radius and thus take smaller gradient updates, thereby initially taking converging steps and converging lower in the loss. Though both methods diverge for α_2 , M-SAM converges remarkably lower in the loss. As M-SAM favors exploiting the current valley rather than exploring the parameter space through larger perturbations, we argue that M-SAM is *conservative* compared to SAM.

Reduced attraction to saddle points. We now study a property of M-SAM that originates from its conservativeness. SAM has the undesirable property of getting attracted to saddle points (Kaddour et al., 2022; Compagnoni et al.,

2023; Kim et al., 2023). Less discussed is the fact that this attraction holds for any type of suboptimal equilibria, meaning that SAM *can* get attracted to maxima as well (Ujváry et al., 2022; Compagnoni et al., 2023). We present an example in Figure 2, where we also show that M-SAM is less prone to get attracted to such equilibria.

We consider the close neighborhood of any parameter configuration θ and define it by a Gaussian with sufficiently low variance σ^2 , i.e. $\tilde{\theta} \sim \mathcal{N}(\theta, \sigma^2 \mathbb{I}_K)$. For simplifying expressions, we reparameterize the neighborhood as $\tilde{\theta} = \theta + \epsilon$ where $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I}_K)$. With the aim of examining attraction to suboptimal equilibria, we describe the general SAM dynamics at $\tilde{\theta}$ using SAM’s gradient flow, eventually giving rise to an ordinary differential equation (ODE) governing the neighborhood:

$$\frac{d\epsilon}{dt} \approx - \underbrace{\left(\nabla^2 \ell(\theta + \delta^*) [\mathbb{I}_K + \tilde{\rho} \cdot \nabla^2 \ell(\theta)] \right)^\top}_{\mathbf{A}_\rho(\theta)} \epsilon. \quad (10)$$

We provide the associated proof in Appendix A. Here, $\tilde{\rho}$ is the general effective perturbation radius and could be replaced by $\tilde{\rho}_{\text{SAM}}$ or $\tilde{\rho}_{\text{M-SAM}}$. In essence, the ODE describes how a random sample within the close neighborhood of θ behaves, yet further analysis of the dynamics described by $\mathbf{A}_\rho(\theta)$ is essential for understanding SAM’s attraction to suboptimal equilibria. First, we remark the minus-sign, revealing that directions in which $\mathbf{A}_\rho(\theta)$ expands, i.e. the directions with positive eigenvalues, instead pull nearby points towards θ . Hence, the point θ is an *attractor* if all eigenvalues of $\mathbf{A}_\rho(\theta)$ are positive. Though the eigenvalues of $\mathbf{A}_\rho(\theta)$ do not have an analytical expression for arbitrary θ , the dynamics simplify remarkably close to equilibria since $\|\nabla \ell(\theta)\|_2 \approx 0$, hence $\delta^* \approx \mathbf{0}$. As argued by Kim et al. (2023) this gives rise to a criteria of stability at equilib-

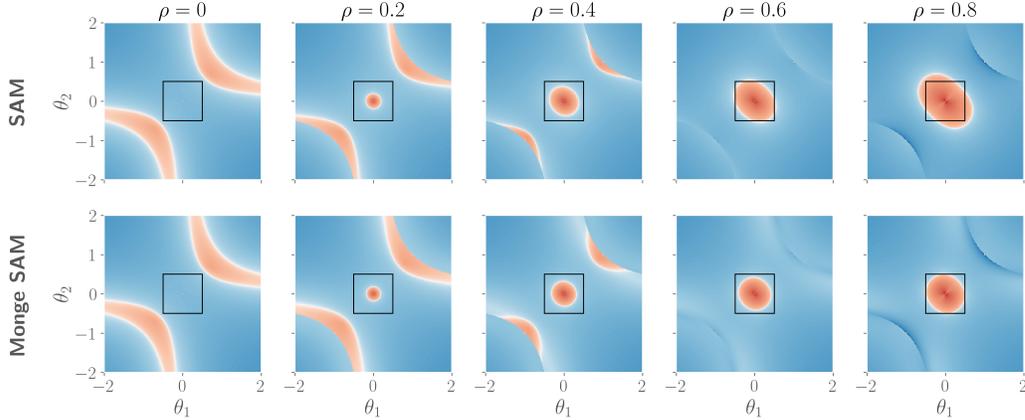


Figure 4. Stability criteria. Intuitively, we highlight the stability of a parameter θ by considering the behavior of a random sample $\tilde{\theta} = \theta + \epsilon$ within the ϵ -neighborhood of θ . We consider the banana-shaped loss, i.e. $\ell(\theta) = (1 - \theta_1\theta_2)^2$, and compute the eigenvalues, $\{\lambda_1, \lambda_2\}$, of $\mathbf{A}_\rho(\theta)$. We show λ_2 across the parameter space for varying perturbation radii for SAM and M-SAM. As λ_1 is always positive, stability of SAM and M-SAM dynamics is satisfied when λ_2 is positive (red regions). The square is used for reference. When $\rho = 0$, i.e. GD behavior, the ϵ -neighborhood of any point in parameter space will converge to the global minima, revealing that GD does not suffer from saddle point attraction if not initialized exactly at the saddle point, $\theta_s = (0, 0)$. Contrarily, θ_s is stable for SAM and M-SAM for higher ρ , even so, M-SAM’s region of attraction is smaller near the saddle but larger near flat global minima at $\theta_{\text{flat}}^* \in \{(-1, -1), (1, 1)\}$, due to its conservative nature.

ria under SAM dynamics depending on the local curvature:

$$\tilde{\rho} > -1/\lambda_i, \quad \forall \lambda_i \in \{\lambda_1, \dots, \lambda_K\}, \quad (11)$$

where $\{\lambda_1, \dots, \lambda_K\}$ are the eigenvalues of the Hessian, $\nabla^2 \ell(\theta)$. The criteria reveals that *any equilibria can be stable* under SAM dynamics if the perturbation is relatively large compared to the local curvature, e.g. if $\tilde{\rho} = 0.35$ and we consider a saddle point or a maximum with eigenvalues of $\{-3, 3\}$ and $\{-3, -3\}$, respectively. This aligns with Figure 2 where the maximum of the sinc-function is a stronger attractor when the perturbation radius is large.

The direct interpretation of M-SAM’s conservativeness, namely that $\tilde{\rho}_{\text{M-SAM}} \leq \tilde{\rho}_{\text{SAM}}$, reveals that the proposed stability criteria is less frequently satisfied for M-SAM since $\tilde{\rho}_{\text{SAM}} \geq \tilde{\rho}_{\text{M-SAM}} > -1/\lambda_i$, and hence *M-SAM is attracted less to suboptimal equilibria than SAM*. We demonstrate empirically this beneficial property of M-SAM in Figure 4 by computing the eigenvalues of $\mathbf{A}_\rho(\theta)$ across the space for the function that we considered in Figure 1. Eigenvalue positivity serves as a criteria for stability and we see that SAM’s region of attraction is larger than M-SAM’s around the saddle point for larger perturbation radii. The behavior is the opposite when close to flat global minima.

3.3. Illustrating Monge SAM: Toy Example in 2D

In Figure 5 we consider the loss $\ell(\theta) = (1 - \theta_1\theta_2)^2$ with banana-shaped minima (as in Figure 1). We initialize SAM, ASAM and M-SAM on a uniform grid $(\theta_1, \theta_2) \in [-3, 3] \times [-3, 3]$ and run the optimizers from each point for 200 steps using a learning rate of 0.01 and a large perturbation radius

of $\rho = 1$. We remark the non-probabilistic formulation, yet we compare the two approaches to Fisher SAM.

All methods are clearly capable of finding the flatter solution when reaching the valleys of global minima; a key property of SAM that motivates its use for fine-tuning tasks and originates from normalizing the gradient in the perturbation step (Dai et al., 2024). Yet, M-SAM reaches the flat global minima more easily than SAM which gets attracted to the saddle point, even when initialized far from it. Fisher SAM appears to be equally robust as M-SAM in this manner and they also behave remarkably similar when considering their marginal distributions of convergence. Only a few of the trajectories close to the center qualitatively differ, which reflects the implications of the ad-hoc approximation of the inverse Fisher metric with a diagonal matrix where M-SAM is equally fast but also captures the correlation.

4. Experiments

We empirically demonstrate how M-SAM behaves in the context of deep neural networks. We rely on mini-batches which introduces stochasticity in the method, whereas the analysis above is based on the deterministic version. We compare with SGD and SAM for understanding M-SAM in high-dimensional settings and include a comparison with Adam on a representational alignment task.

4.1. Fine-tuning from a Bad Global Minimum

We consider a pre-trained ResNet-18 for classifying images of CIFAR10, converged to a bad global minima (Liu

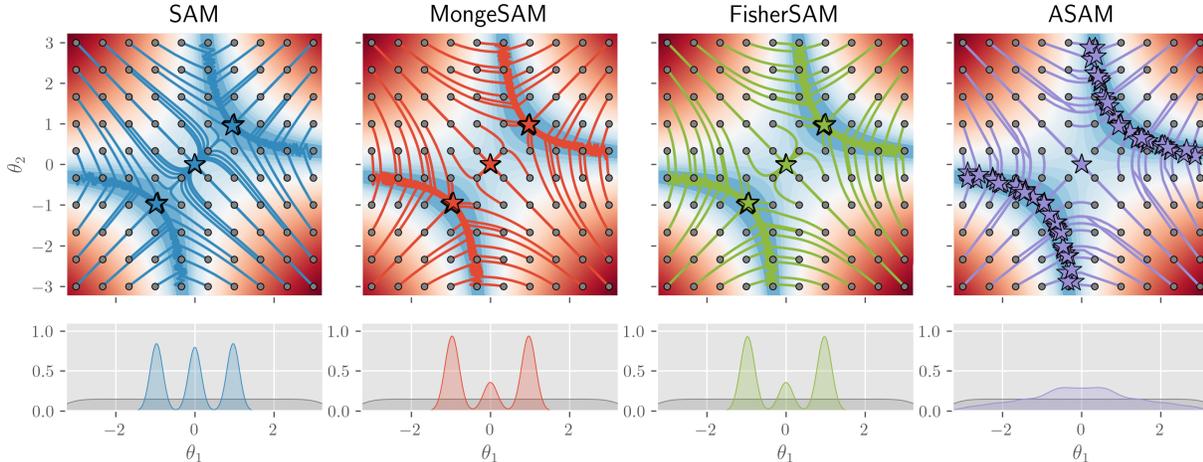


Figure 5. Comparing sharpness-aware minimizers. Trajectories of SAM, M-SAM, Fisher SAM and ASAM in the banana-shaped loss given by $\ell(\theta) = (1 - \theta_1\theta_2)^2$. Each optimizer is initialized on the points of a 10×10 grid and runs to convergence. The marginal distributions of convergence locations summarize the final estimates. While all methods can explore the valley of minima and locate the flattest solution, SAM exhibits a notable tendency to converge to the saddle point. M-SAM and Fisher SAM show qualitatively similar behavior, with only minor differences when initialized near the saddle. We note that this behavior is due to the diagonal approximation of the Fisher metric. We compare with ASAM for completeness; it mostly avoids the saddle point but converges anywhere on the minima range and does generally not find one of the flattest solution.

et al., 2020) that perfectly fits the training data while obtaining 48% generalization accuracy. We fine-tune with SGD, SAM and M-SAM using a small learning rate, i.e. $\alpha = 0.001$, with the aim of restricting the optimizers to take steps inside the valley containing the bad global minimum, similarly to Dai et al. (2024). We use a fixed batch size of $B = 128$ and train for 500 epochs using a NLL loss on log-softmax outputs of the network. We experiment with $\rho \in \{0.005, 0.01, 0.03, 0.05\}$. Even if some configurations do not converge, we deem the training horizon sufficiently long for examining the dynamics of the methods. The results are presented in Figure 6 with uncertainty estimates empirically computed as the standard error of the mean (SEM) from $N = 5$ pseudo-random repetitions.

For lower perturbation radii $\rho \in \{0.005, 0.01, 0.03\}$, SAM and M-SAM behave similarly as in the 2D toy setting; they are capable of moving along the valley of bad global minima, eventually converging to a solution with increased generalization performance. Specifically, searching wider with $\rho = 0.03$ allows for improving approx. 6% compared to fine-tuning with SGD. Though increasing the perturbation size to $\rho = 0.05$ increases the model’s generalization capabilities even further it also reveals the main controversial property of M-SAM, namely the difficulty to escape the valley of bad global minima. Contrarily, SAM with large perturbations (i.e. $\rho = 0.05$) quickly escapes the initial valley and finds a region of the parameter space where the generalization performance significantly improves.

We emphasize that this behavior which leads to better per-

formance, in a different scenario might come at the cost of forgetting what was learned during pre-training, potentially making M-SAM a better choice for fine-tuning. As we showed above, large perturbation radii increase SAM’s attraction to suboptimal equilibria (Section 3.2); so applying SAM with uncalibrated perturbation parameter ρ to fine-tune pre-trained models that generalize well, potentially in sharp parameter regions, may ultimately result in convergence at another suboptimal location e.g., saddle point. In fact, we provide such an experiment below.

4.2. How does M-SAM affect cross-modal alignment?

A recent study (Huh et al., 2024) suggested that the alignment of the latent representation space of multi-modal models improves with the capacity of the domain-specific encoders. We examine whether fine-tuning with a sharpness-aware minimizer can improve representational alignment by finding better minima, potentially leading to more robust representations and improved alignment, without requiring increasing the model capacity.

We fine-tune a transformer-based CLIP model (Radford et al., 2021) on the Wikimedia version of the WIT dataset (available at huggingface.co) with SGD, Adam, SAM and M-SAM using CLIP loss. We used a batch size of 128 and train to convergence. We run a grid-search for SAM and M-SAM with combinations of learning rates, α , and perturbation radii, ρ , from $\{1 \cdot 10^{-5}, 5 \cdot 10^{-5}, 1 \cdot 10^{-4}\}$ and $\{0.01, 0.02, 0.03, 0.04\}$, respectively. For SAM we also consider $\rho = 0.005$. For SGD and SAM we examine learn-

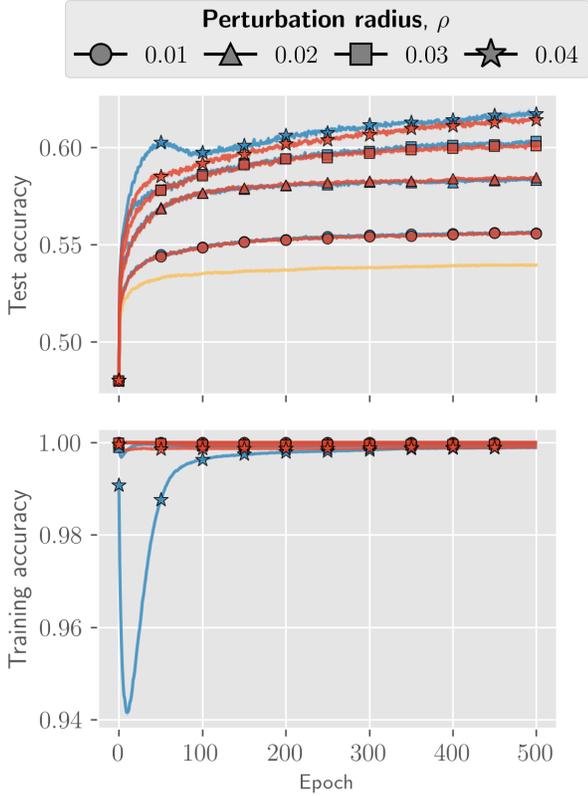


Figure 6. Fine-tuning a bad ResNet-18. We initialize an adversarially trained ResNet-18 (Liu et al., 2020) using SGD (—), SAM (—) and M-SAM (—) with a learning rate set to 0.001. We examine SAM and M-SAM behavior for varying perturbation radii and see that M-SAM’s conservativeness biases the model to exploit the current valley of global minima rather than escaping and exploring other regions of the parameter space. We provide uncertainty estimates as the standard error of the mean computed from $N = 5$ pseudo-random repetitions. Remark that most models perfectly fit the training data, retaining a training accuracy of 100%.

ing rates in the range $[10^{-6}, 10^{-2}]$, yet most settings converge slowly to worse solutions than SAM-based methods or immediately diverge. We evaluate the generalization capabilities of the pre-trained and fine-tuned models on a subset of the MS-COCO Captions dataset (Lin et al., 2014), restricted to having one caption per image. We measure alignment with the mutual k NN similarity score, \mathcal{S}_{kNN} , (Huh et al., 2024) by considering the fraction of $k = 8$ neighbors shared for the latent text and image representation for an image-text input pair; this captures whether conceptual regions form in CLIP space. We report the optimal alignment in Table 1 and a comparison of CLIP losses at convergence for SAM and M-SAM in Table 2 when fixing the perturbation radius to $\rho = 0.01$. The remaining results are provided in Appendix B. We also include a visualization of the latent representations space in Figure 7 for qualitatively assessing the impact of fine-tuning with SAM and M-SAM.

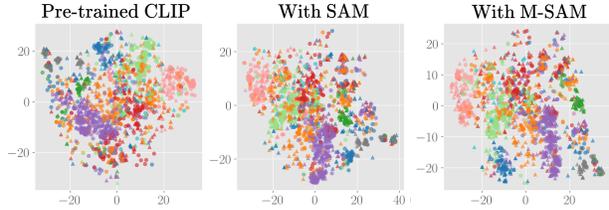


Figure 7. Visualizing CLIP space. We stack the latent text and image representations (MS-COCO) and project them with PCA to a lower-dimensional subspace. The main variation in the stacked representation space originates from the modality. So we remove the information of the first principal component direction before running tSNE on the modified reconstructed latent representations. We label each image-text pair by classifying the images using a ResNet-50, pre-trained on ImageNet labels and color them according to their superclass label (e.g. *Vehicle*, *Clothing* or *Food*), obtained by backtracking the ImageNet graph. We restrict ourselves to consider image-text pairs from the 12 major superclasses.

Table 1. Alignment scores of fine-tuned CLIP models on MS-COCO. Optimal hyperparameters are shown in Appendix B.

	None	SGD	Adam	SAM	M-SAM
\mathcal{S}_{kNN}	0.351	0.387	0.388	0.405	0.446

Table 2. Fine-tuned CLIP loss after convergence for varying learning rates with $\rho = 0.01$. Remarkable divergence is shown in red.

	$1 \cdot 10^{-5}$	$5 \cdot 10^{-5}$	$1 \cdot 10^{-4}$
SAM	0.607	4.852	4.852
M-SAM	0.584	0.584	0.584

Though SGD and Adam improves representational alignment with approx. 3-4% compared to the pre-trained model - revealing the benefit of fine-tuning for improving downstream performance - seeking flatter minima with SAM and M-SAM gives even better conceptual alignment, with M-SAM finding a remarkably better solution than SAM. The fact that M-SAM learns a conceptually better aligned representation space than SAM relates to M-SAM’s robustness to hyperparameter tuning; increasing the perturbation radius for M-SAM generally enhances representational alignment, whereas SAM easily becomes unstable and take diverging steps, supposedly converging to a saddle point or local maxima. Thus, M-SAM’s conservativeness appears to also practically introduce robustness to hyperparameter tuning, potentially allowing for improved downstream performance.

5. Related Work

Extensions of SAM. As mentioned, SAM contains an inherent scale-dependency problem (Dinh et al., 2017). A

widely applied approach (Kwon et al., 2021) - namely, adaptive SAM (ASAM) - exploits a scale-invariant sharpness definition and modifies SAM by computing the worst-case perturbation in a normalized parameter space, typically by element-wise parameter scaling. Though improving on SAM on a variety of tasks including image classification, robustness to label noise and machine translation, the normalization operator is defined in an ad-hoc manner. This - along with the fact that the parameter space of deep neural networks is typically a statistical manifold that is not properly captured by Euclidean geometry - motivated the development of Fisher SAM (Kim et al., 2022) that obtains slightly improved generalization performance on CIFAR-10 and CIFAR-100 than SAM and ASAM for a range of vision backbones. Recent advances include Riemannian SAM (Yun & Yang, 2024); a general framework containing Fisher SAM that applies SAM on Riemannian manifolds by exploiting the tangent space and exponential maps. Applying the method requires predefining the manifold on which to optimize. They provide convergence analyses for SAM under the general formulation on Riemannian manifolds.

Other approaches to find flat minima. Stochastic weight averaging (SWA) (Izmailov et al., 2018) exploits averaging of weights over iterations of SGD; SWA is shown to find wider solutions than SGD and empirically improved generalization compared to SGD for a variety of tasks. Meanwhile, SAM can find flatter optima than SWA in terms of eigenvalues of the Hessian (Kaddour et al., 2022), yet, the better approach in terms of generalization performance is data- and task-specific with SWA e.g. working better for graph representation learning tasks than SAM. Other approaches include adjusting the gradient update by adding information from a random perturbation when the gradient norm is sufficiently low (Ahn et al., 2023) or applying SGD with dominant noise (Keskar et al., 2016; Jastrzebski et al., 2017; Xing et al., 2018; Zhu et al., 2018; Smith et al., 2020; Zhang et al., 2021), e.g. through large step sizes or small batches.

The Monge metric. In recent years, the Monge metric has been used for various applications due to its simplistic form. It was comprehensively studied in context of geometric Markov Chain Monte Carlo (MCMC) sampling (Hartmann et al., 2022) to replace the Fisher information matrix (FIM) for improving efficiency of geometric MCMC sampling. The study revealed that the metric has an inherent notion of curvature. A different study (Bergamin et al., 2024) proposed a Riemannian Laplace approximation for Bayesian neural networks (BNN) using the Monge metric for better capturing the underlying geometry of the intractable posterior. The method worked well on a variety of problems, yet the conservative nature of the metric implies suboptimal samples (Yu et al., 2023). Instead Yu et al. (2023) propose replacing the Monge metric with a metric based on the

FIM, resulting in improved sample quality, which increases substantially the computational complexity.

6. Conclusion

In this paper, we proposed Monge SAM (M-SAM), a novel reparameterization-invariant approach to sharpness-aware minimization. M-SAM leverages the Monge metric, a Riemannian metric in the parameter space naturally induced by the training loss surface. The method works under any modeling choice, relying only on a smoothness assumption of the loss function.

Intuitively, the proposed M-SAM searches for the worst-case adversarial perturbation according to the Monge metric, resulting in “conservative” perturbations bounded by SAM’s perturbations. This makes M-SAM behave as a mixture of SAM and GD, with the specific implications that we analyzed and observed empirically in the context of deep neural networks.

The experiments revealed that while SAM escapes valleys more effectively to find flatter solutions with potentially better generalization, M-SAM is more robust to hyperparameter selection, namely learning rate and perturbation radii, which improves the algorithmic stability. For fine-tuning tasks, particularly in improving multi-modal representational alignment, M-SAM converged to solutions that improve the representational alignment remarkably while SAM exhibited instability, likely due to SAM’s tendency to converge to suboptimal equilibria like saddle points or maxima. We further provided theoretical evidence supporting M-SAM’s resilience against getting attracted to suboptimal equilibria, compared to SAM.

Limitations and future work. Though the conservativeness of M-SAM guarantees properties like robustness and better resilience to saddle point attraction than SAM, it potentially prevents M-SAM from escaping suboptimal valleys for finding globally flatter regions of the parameter space.

Similar to SAM, M-SAM relies on a Taylor expansion to approximate adversarial perturbations, imposing ellipsoidal constraints rather than respecting the true geometric structure. Future work could investigate averaging the Monge metric over iterations or sampling it from the local neighborhood to better capture the local geometry. Another interesting direction is towards building a deeper understanding of the connection between the Monge and Fisher metrics, as previously hinted in related work (Yu et al., 2023).

Impact Statement

This paper introduces a novel optimizer designed to improve the generalization performance of machine learning models while considering the loss geometry. Our work advances optimization techniques in machine learning and has broad applicability. While the method requires more compute than standard gradient descent schemes, it is more robust to hyperparameter settings compared to SAM and might require less hyperparameter tuning when applied in practice. We do not identify any direct negative societal impacts of this contribution.

References

- Ahn, K., Jadbabaie, A., and Sra, S. How to escape sharp minima. *arXiv preprint arXiv:2305.15659*, 2023.
- Bergamin, F., Moreno-Muñoz, P., Hauberg, S., and Arvanitidis, G. Riemannian laplace approximations for bayesian neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019.
- Compagnoni, E. M., Biggio, L., Orvieto, A., Proske, F. N., Kersting, H., and Lucchi, A. An sde for modeling sam: Theory and insights. In *International Conference on Machine Learning*, pp. 25209–25253. PMLR, 2023.
- Dai, Y., Ahn, K., and Sra, S. The crucial role of normalization in sharpness-aware minimization. *Advances in Neural Information Processing Systems*, 36, 2024.
- Dinh, L., Pascanu, R., Bengio, S., and Bengio, Y. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, pp. 1019–1028. PMLR, 2017.
- Dziugaite, G. K. and 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.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Hartmann, M., Girolami, M., and Klami, A. Lagrangian manifold monte carlo on monge patches. In *International Conference on Artificial Intelligence and Statistics*, pp. 4764–4781. PMLR, 2022.
- Hochreiter, S. and Schmidhuber, J. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- Huh, M., Cheung, B., Wang, T., and Isola, P. The platonic representation hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
- Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. Finding flatter minima with sgd. 2018.
- Jiang, Y., Neyshabur, B., Mobahi, H., Krishnan, D., and Bengio, S. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- Kaddour, J., Liu, L., Silva, R., and Kusner, M. J. When do flat minima optimizers work? *Advances in Neural Information Processing Systems*, 35:16577–16595, 2022.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Kim, H., Park, J., Choi, Y., and Lee, J. Stability analysis of sharpness-aware minimization. *arXiv preprint arXiv:2301.06308*, 2023.
- Kim, M., Li, D., Hu, S. X., and Hospedales, T. Fisher sam: Information geometry and sharpness aware minimisation. In *International Conference on Machine Learning*, pp. 11148–11161. PMLR, 2022.
- Kristiadi, A., Dangel, F., and Hennig, P. The geometry of neural nets’ parameter spaces under reparameterization. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kwon, J., Kim, J., Park, H., and Choi, I. K. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning*, pp. 5905–5914. PMLR, 2021.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.

- Liu, H., Xie, S. M., Li, Z., and Ma, T. Same pre-training loss, better downstream: Implicit bias matters for language models. In *International Conference on Machine Learning*, pp. 22188–22214. PMLR, 2023.
- Liu, S., Papailiopoulos, D., and Achlioptas, D. Bad global minima exist and sgd can reach them. *Advances in Neural Information Processing Systems*, 33:8543–8552, 2020.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Smith, S., Elsen, E., and De, S. On the generalization benefit of noise in stochastic gradient descent. In *International Conference on Machine Learning*, pp. 9058–9067. PMLR, 2020.
- Ujváry, S., Telek, Z., Kerekes, A., Mészáros, A., and Huszár, F. Rethinking sharpness-aware minimization as variational inference. *arXiv preprint arXiv:2210.10452*, 2022.
- Xing, C., Arpit, D., Tsirigotis, C., and Bengio, Y. A walk with sgd. *arXiv preprint arXiv:1802.08770*, 2018.
- Yu, H., Hartmann, M., Williams, B., Girolami, M., and Klami, A. Riemannian laplace approximation with the fisher metric. *arXiv preprint arXiv:2311.02766*, 2023.
- Yun, J. and Yang, E. Riemannian sam: Sharpness-aware minimization on riemannian manifolds. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Zhu, Z., Wu, J., Yu, B., Wu, L., and Ma, J. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. *arXiv preprint arXiv:1803.00195*, 2018.

A. Proof of ODE describing local SAM dynamics

We extend the discrete SAM update to continuous time, hence establishing the SAM gradient flow:

$$\frac{d\boldsymbol{\theta}}{dt} = -f(z(\boldsymbol{\theta}))$$

where $f(z(\boldsymbol{\theta})) = \nabla \ell(z(\boldsymbol{\theta}))$ and $z(\boldsymbol{\theta}) = \boldsymbol{\theta} + \tilde{\rho} \nabla \ell(\boldsymbol{\theta})$ and $\tilde{\rho}$ is the effective perturbation radius of either SAM or M-SAM. Next, we examine the dynamics in an ϵ -close neighborhood to $\boldsymbol{\theta}$ where ϵ is a random vector drawn from $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ with sufficiently small variance. We approach this by a Taylor expansion around $\boldsymbol{\theta}$:

$$\begin{aligned} \frac{d\boldsymbol{\theta}}{dt} &\approx \frac{d(\boldsymbol{\theta} + \epsilon)}{dt} = -f(z(\boldsymbol{\theta} + \epsilon)) \approx -\left[f(z(\boldsymbol{\theta})) + \nabla f(z(\boldsymbol{\theta}))^\top \epsilon \right] \\ &= -f(z(\boldsymbol{\theta})) - \nabla f(z(\boldsymbol{\theta}))^\top \epsilon \\ &= \cancel{\frac{d\boldsymbol{\theta}}{dt}} - \nabla f(z(\boldsymbol{\theta}))^\top \epsilon \end{aligned}$$

The dynamics at $\boldsymbol{\theta}$ cancel out and the ϵ -neighborhood is described by an ODE. Inserting f and using the chain rule gives:

$$\begin{aligned} \frac{d\epsilon}{dt} &\approx -\nabla [\nabla \ell(z(\boldsymbol{\theta}))]^\top \epsilon \\ &= -\frac{\partial}{\partial \boldsymbol{\theta}} [\nabla \ell(z(\boldsymbol{\theta}))]^\top \epsilon \\ &= -\left(\frac{\partial}{\partial z(\boldsymbol{\theta})} \nabla \ell(z(\boldsymbol{\theta})) \cdot \frac{\partial z(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^\top \epsilon \\ &= -\left(\frac{\partial}{\partial z(\boldsymbol{\theta})} \nabla \ell(z(\boldsymbol{\theta})) \cdot \left[\frac{\partial}{\partial \boldsymbol{\theta}} \boldsymbol{\theta} + \tilde{\rho} \frac{\partial}{\partial \boldsymbol{\theta}} \nabla \ell(\boldsymbol{\theta}) \right] \right)^\top \epsilon \\ &= -\underbrace{\left(\nabla^2 \ell(z(\boldsymbol{\theta})) \cdot [\mathbf{I}_d + \tilde{\rho} \nabla^2 \ell(\boldsymbol{\theta})] \right)^\top}_{\mathbf{A}_\rho(\boldsymbol{\theta})} \epsilon \end{aligned}$$

where $\nabla^2 \ell(\boldsymbol{\theta})$ is the Hessian at $\boldsymbol{\theta}$. Thus, SAM dynamics are governed by the local curvature at $\boldsymbol{\theta}$ along with the curvature at the perturbed parameter set $z(\boldsymbol{\theta}) = \boldsymbol{\theta} + \tilde{\rho} \nabla \ell(\boldsymbol{\theta})$ when close to $\boldsymbol{\theta}$. Examining eigenvalues of $\mathbf{A}_\rho(\boldsymbol{\theta})$ allows for addressing stability of SAM across the parameter space. Due to the minus sign, cases where the eigenvalues of $\mathbf{A}_\rho(\boldsymbol{\theta})$ are positive gives decaying dynamics, meaning $\epsilon \rightarrow 0$, thus nearby points approach $\boldsymbol{\theta}$ and we call $\boldsymbol{\theta}$ an *attractor* under SAM dynamics.

What Kim et al. (2023) do, is to consider stability at equilibria points of the loss landscape, i.e. $\boldsymbol{\theta}^*$ where $\nabla \ell(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} = 0$. In such settings the SAM perturbation disappears as $z(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} = \boldsymbol{\theta}$ which results in

$$\begin{aligned} \mathbf{A}_\rho(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} &= -\nabla^2 \ell(\boldsymbol{\theta}^*) \cdot [\mathbf{I}_d + \tilde{\rho} \nabla^2 \ell(\boldsymbol{\theta}^*)] \\ &= -\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top \cdot [\mathbf{I} + \tilde{\rho}\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top] \\ &= -\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top + \tilde{\rho}\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top \\ &= -\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top + \tilde{\rho}\mathbf{Q}\boldsymbol{\Lambda}^2\mathbf{Q}^\top \\ &= -\mathbf{Q}[\boldsymbol{\Lambda} + \tilde{\rho}\boldsymbol{\Lambda}^2]\mathbf{Q}^\top \end{aligned}$$

where the Hessian matrix, $\nabla^2 \ell(\boldsymbol{\theta}^*)$, is real and symmetric and thus factorizes through the eigendecomposition to the form $H_\ell(\boldsymbol{\theta}^*) = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top$, where \mathbf{Q} is an orthonormal matrix, i.e. $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$ contains the eigenvalues. Specifically for equilibria, Kim et al. (2023) determine a stability criteria by considering the eigenvalues of the Hessian. This means that all eigenvalues should satisfy:

$$-(\lambda_i + \tilde{\rho}\lambda_i^2) = -\lambda_i - \tilde{\rho}\lambda_i^2 < 0, \quad \forall \lambda_i \in \{\lambda_1, \dots, \lambda_d\}$$

Since $\lambda_i^2 \geq 0$, isolating for $\tilde{\rho}$ gives:

$$\tilde{\rho} > -\frac{1}{\lambda_i}, \quad \forall \lambda_i \in \{\lambda_1, \dots, \lambda_d\}$$

B. Multi-modal representational alignment - full grid search

Table 3. **CLIP performance after fine-tuning.** We evaluate pre-trained CLIP on the COCO Captions dataset as well as models fine-tuned to the Wiki dataset. We consider SGD, Adam, SAM and M-SAM with varying (α, ρ) -settings where α is the learning rate and ρ is the perturbation radius. We present the loss and alignment scores at their optimal iterations, T_{loss}^* and $T_{\text{alignment}}^*$, and mark the optimal loss values for each optimizer type in violet. We highlight diverging behavior as red cells with * denoting divergence behavior from the first fine-tuning step. Green cells represent more than 5% improvement in terms of representational alignment, compared to the pre-trained CLIP model and the overall optimal performance is marked in bold text.

Optimizer	Params		Alignment (\uparrow)	Loss (\downarrow)	
	α	ρ	$T = T_{\text{alignment}}^*$	$T = 2000$	$T = T_{\text{loss}}^*$
None	-	-	*0.351	-	*0.699
SGD	$1 \cdot 10^{-4}$	-	0.376	0.596	0.591
	$1 \cdot 10^{-3}$	-	0.387	0.663	0.634
Adam	$1 \cdot 10^{-6}$	-	0.376	0.631	0.595
	$5 \cdot 10^{-6}$	-	0.388	0.600	0.600
SAM	$1 \cdot 10^{-5}$	0.005	0.397	0.597	0.584
		0.01	0.405	0.623	0.607
		0.02	*0.351	4.849	*0.699
		0.03	*0.351	5.077	0.697
		0.04	*0.351	5.503	*0.699
	$5 \cdot 10^{-5}$	0.005	0.401	0.593	0.591
		0.01	*0.351	4.852	*0.699
		0.02	*0.351	4.852	*0.699
		0.03	*0.351	5.026	*0.699
		0.04	*0.351	4.882	*0.699
	$1 \cdot 10^{-4}$	0.005	0.403	0.602	0.594
		0.01	*0.351	4.852	*0.699
		0.02	*0.351	4.852	*0.699
		0.03	*0.351	4.870	*0.699
		0.04	*0.351	4.866	*0.699
	M-SAM	$1 \cdot 10^{-5}$	0.01	0.388	0.586
0.02			0.411	0.594	0.578
0.03			0.415	0.588	0.576
0.04			0.420	0.620	0.604
$5 \cdot 10^{-5}$		0.01	0.404	0.591	0.584
		0.02	0.417	0.615	0.606
		0.03	0.433	0.602	0.587
		0.04	0.446	0.620	0.615
$1 \cdot 10^{-4}$		0.01	0.405	0.602	0.584
		0.02	0.422	0.628	0.608
		0.03	*0.351	4.852	*0.699
		0.04	*0.351	4.852	*0.699