

Penalizing Gradient Norm for Efficiently Improving Generalization in Deep Learning

Yang Zhao¹ & Hao Zhang¹ & Xiuyuan Hu¹

¹ Department of Electronic Engineering, Tsinghua University

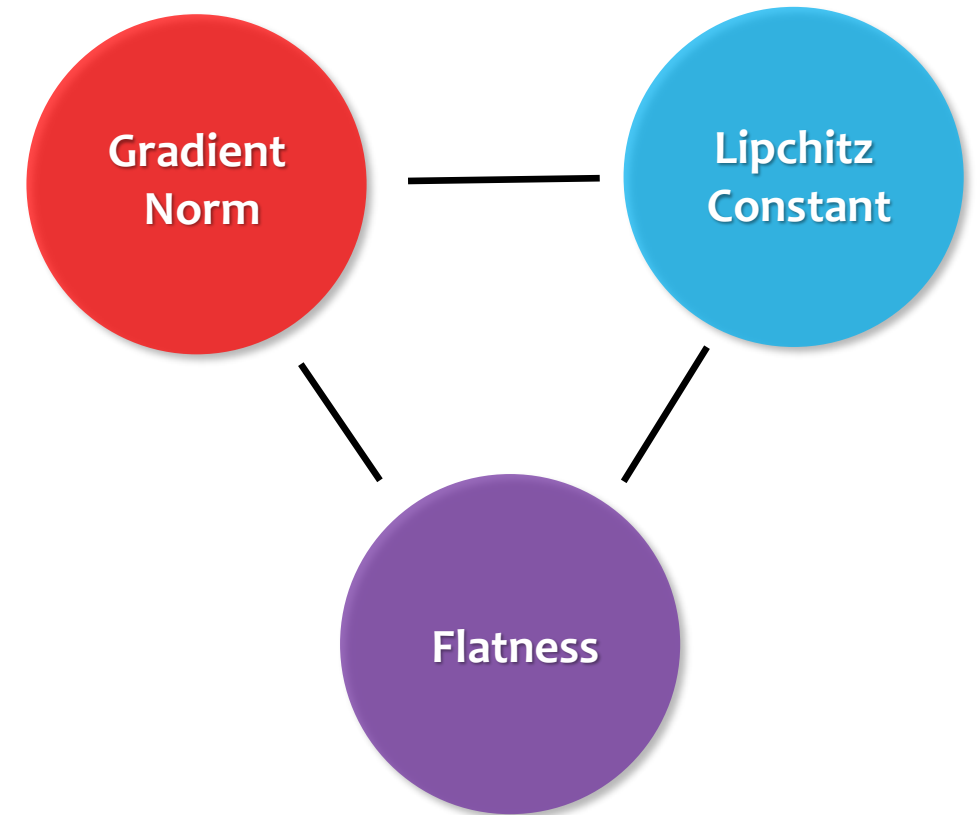
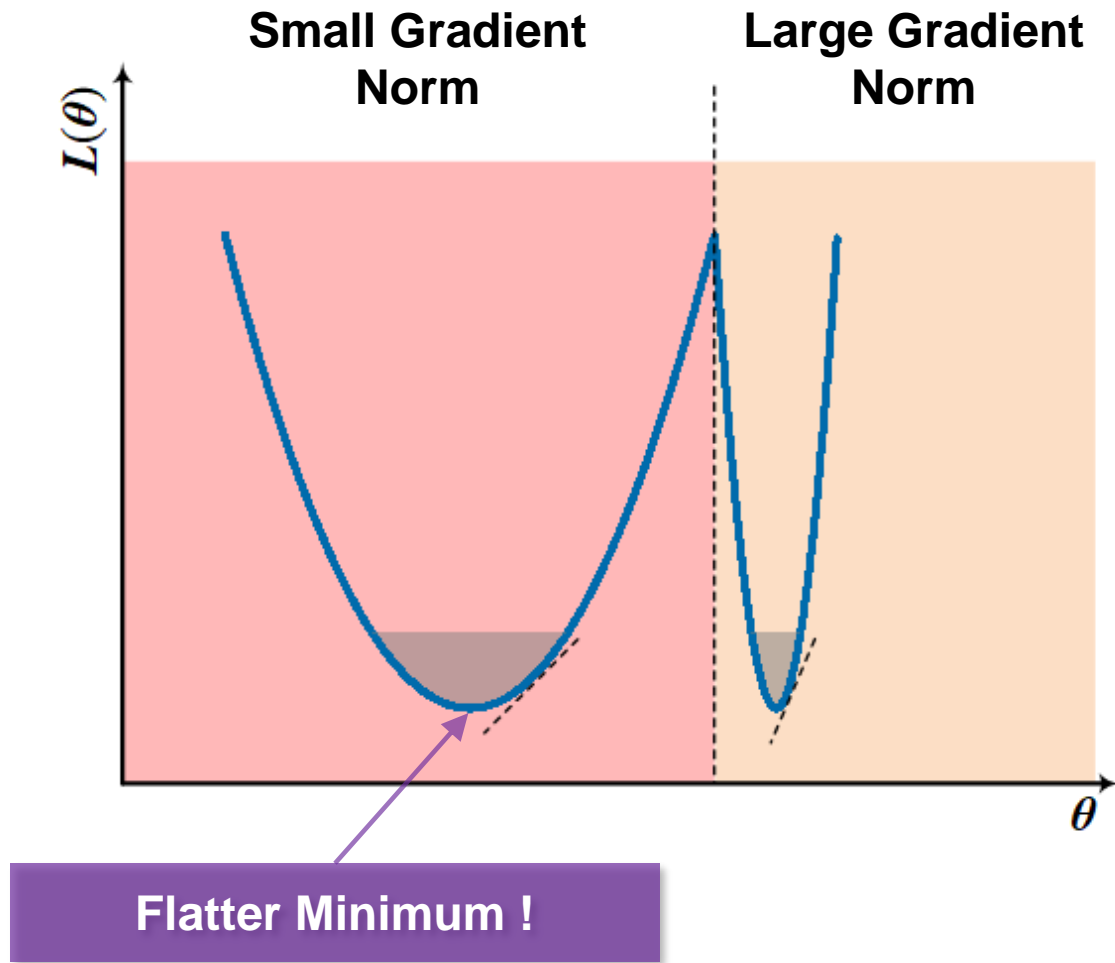
- Basic setting

$$L(\boldsymbol{\theta}) = \underbrace{L_{\mathcal{S}}(\boldsymbol{\theta})}_{\text{Empirical loss}} + \underbrace{\lambda \cdot \|\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})\|_2}_{\text{Penalty coefficient} \cdot \text{Gradient norm of empirical loss}}$$

Diagram illustrating the components of the loss function $L(\boldsymbol{\theta})$:

- $L_{\mathcal{S}}(\boldsymbol{\theta})$ is the **Empirical loss**.
- λ is the **Penalty coefficient**.
- $\|\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})\|_2$ is the **Gradient norm of empirical loss**.

- Motivation



- Computing Gradient

$$L(\boldsymbol{\theta}) = L_{\mathcal{S}}(\boldsymbol{\theta}) + \lambda \cdot ||\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})||_2$$

$\nabla_{\boldsymbol{\theta}}$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta}) + \lambda \nabla_{\boldsymbol{\theta}} ||\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})||_2$$

Chain Rule

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta}) + \lambda \cdot \nabla_{\boldsymbol{\theta}}^2 L_{\mathcal{S}}(\boldsymbol{\theta}) \frac{\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})}{||\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})||}$$

Hessian Computation !

- **Approximation of** $\nabla_{\theta}^2 L_{\mathcal{S}}(\theta) \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|}$

Taylor expansion :

$$\nabla_{\theta} L_{\mathcal{S}}(\theta + \Delta\theta) = \nabla_{\theta} L_{\mathcal{S}}(\theta) + \mathbf{H} \Delta\theta + \mathcal{O}(\|\Delta\theta\|^2)$$

Make $\Delta\theta = r\mathbf{v}$

$$\mathbf{H}\mathbf{v} = \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta + r\mathbf{v}) - \nabla_{\theta} L_{\mathcal{S}}(\theta)}{r} + \mathcal{O}(r)$$

Make $\mathbf{v} = \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|}$

$$\nabla_{\theta}^2 L_{\mathcal{S}}(\theta) \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|} \approx \frac{\nabla_{\theta} L(\theta + r \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|}) - \nabla_{\theta} L(\theta)}{r}$$

- Practical Implementation

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta}) + \frac{\lambda}{r} \cdot (\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta} + r \frac{\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})}{\|\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})\|}) - \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})) \\ &= (1 - \alpha) \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta}) + \alpha \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta} + r \frac{\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})}{\|\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})\|}) \\ &\approx (1 - \alpha) \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta}) + \alpha \nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta} + r \frac{\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})}{\|\nabla_{\boldsymbol{\theta}} L_{\mathcal{S}}(\boldsymbol{\theta})\|}}\end{aligned}$$

Note :

1. The parameter λ denotes the penalty coefficient, controlling the degree of penalty. We should set it properly during training.
2. The parameter r is used for approximating the Hessian vector multiplication. It should be set smaller enough to ignore the big O item safely. But on the other side, as r becomes smaller, the other term $r\mathbf{v}$ will approach zero, causing a loss of precision on \mathbf{v} in approximating $H\mathbf{v}$.

- Practical Implementation

Algorithm 1:

Full procedures of our optimization scheme. This algorithm only shows gradient norm penalty with SGD as its base optimizer. For other optimizers or update strategies (such as Adam), one could add specific operations before step 8.

Algorithm 1 Optimization Scheme of Penalizing Gradient Norm

Input: Training set $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^N$; loss function $L(\cdot)$; batch size B ; learning rate η ; total step T ; balance coefficient α ; approximation scalar r .

Parameter: Model parameters θ

Output: Optimized weight $\hat{\theta}$

- 1: Parameter initialization θ_0 .
- 2: **for** step $t = 1$ **to** T **do**
- 3: Get batch data pairs $\mathcal{B} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^B$ sampled from training set \mathcal{S} .

| | |
|--------|---|
| Step 1 | 4: Calculate the gradient $\mathbf{g}_1 = \nabla_{\theta} L_{\mathcal{S}}(\theta)$ based on the batch samples. |
| Step 2 | 5: Add $r \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\ \nabla_{\theta} L_{\mathcal{S}}(\theta)\ }$ on the current parameter θ_t , which makes $\theta_t' = \theta_t + r \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\ \nabla_{\theta} L_{\mathcal{S}}(\theta)\ }$. |
| | 6: Calculate the gradient $\mathbf{g}_2 = \nabla_{\theta} L_{\mathcal{S}}(\theta)$ at $\theta = \theta_t'$. |
| Step 3 | 7: Calculate the final gradient $\mathbf{g} = (1 - \alpha)\mathbf{g}_1 + \alpha\mathbf{g}_2$. |
| | 8: (SGD optimizer) Update parameter with final gradient, $\theta_{k+1} = \theta_k - \eta \cdot \mathbf{g}$. |

9: **end for**

10: **return** Final optimization $\hat{\theta}$.

- Connection with SAM

Gradient Norm Penalty

- Optimization target :

$$\min_{\theta} L_{\mathcal{S}}(\theta) + \lambda \cdot \|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|_2$$

- Optimization solution :

$$\begin{aligned} \nabla_{\theta} L(\theta) &= (1 - \alpha) \nabla_{\theta} L_{\mathcal{S}}(\theta) \\ &\quad + \alpha \nabla_{\theta} L_{\mathcal{S}}(\theta) \Big|_{\theta = \theta + r \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|}} \end{aligned}$$

Sharpness-Aware Minimization^[1]

- Optimization target :

$$\min_{\theta} \max_{\|\epsilon\|_2 \leq r} L_{\mathcal{S}}(\theta + \epsilon)$$

- Optimization solution :

$$\begin{aligned} \nabla_{\theta} L(\theta) &= 0 \cdot \nabla_{\theta} L_{\mathcal{S}}(\theta) \\ &\quad + 1 \cdot \nabla_{\theta} L_{\mathcal{S}}(\theta) \Big|_{\theta = \theta + r \frac{\nabla_{\theta} L_{\mathcal{S}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{S}}(\theta)\|}} \end{aligned}$$

Apparently, SAM is a special implementation of gradient norm penalty, where $\alpha = 1$. In other words, SAM will create connections between the two independent parameters, where the coefficient λ will be always set equal to r . This limits our training and models could not get the best performance for $\alpha = 1$ in gradient norm penalty.

- Training CNN Models on Cifar

Table :

Testing error rate of CNN models on Cifar10 and Cifar100 when training with the standard SGD, the SAM and our scheme, respectively.

It could be found that SAM could outperform the standard scheme, but could not be the best implementation in our scheme..

| | Cifar10 | | Cifar100 | |
|-----------------------------|---------------------------------|---------------------------------|----------------------------------|----------------------------------|
| VGG16 | Basic | Cutout | Basic | Cutout |
| Standard | 7.07 | 5.31 | 28.78 | 26.98 |
| SAM | 6.91 | 6.17 | 28.62 | 27.13 |
| Ours | 6.72 | 5.19 | 28.48 | 27.07 |
| VGG16-BN | Basic | Cutout | Basic | Cutout |
| Standard | 5.74 \pm 0.09 | 4.39 \pm 0.07 | 25.22 \pm 0.31 | 24.69 \pm 0.25 |
| SAM | 5.24 \pm 0.08 | 4.16 \pm 0.11 | 24.23 \pm 0.29 | 23.35 \pm 0.33 |
| Ours | 4.88\pm0.12 | 4.02\pm0.08 | 24.04\pm0.18 | 23.07\pm0.26 |
| WideResNet-28-10 | Basic | Cutout | Basic | Cutout |
| Standard | 3.53 \pm 0.10 | 2.81 \pm 0.07 | 18.99 \pm 0.12 | 16.92 \pm 0.10 |
| SAM | 2.78 \pm 0.07 | 2.43 \pm 0.13 | 16.53 \pm 0.13 | 14.87 \pm 0.16 |
| Ours | 2.52\pm0.09 | 2.16\pm0.11 | 16.02\pm0.19 | 14.28\pm0.16 |
| WideResNet-SS 2 \times 96 | Basic | Cutout | Basic | Cutout |
| Standard | 2.82 \pm 0.05 | 2.39 \pm 0.06 | 17.19 \pm 0.19 | 15.85 \pm 0.14 |
| SAM | 2.37 \pm 0.09 | 2.11 \pm 0.13 | 15.22 \pm 0.19 | 14.32 \pm 0.15 |
| Ours | 2.28\pm0.13 | 2.01\pm0.10 | 14.93\pm0.10 | 14.03\pm0.17 |
| PyramidNet-SD | Auto Aug + Cutmix | | Auto Aug + Cutmix | |
| Standard | 1.66 \pm 0.11 | | 10.83 \pm 0.14 | |
| SAM | 1.41 \pm 0.08 | | 10.33 \pm 0.13 | |
| Ours | 1.30\pm0.07 | | 10.12\pm0.17 | |

- Training ViT Models on Cifar

Table :

Testing error rate of ViT models on Cifar10 and Cifar100 when training with the standard SGD, the SAM and our scheme, respectively.

It could be found that SAM could outperform the standard scheme, but again it could not be the best implementation in our scheme..

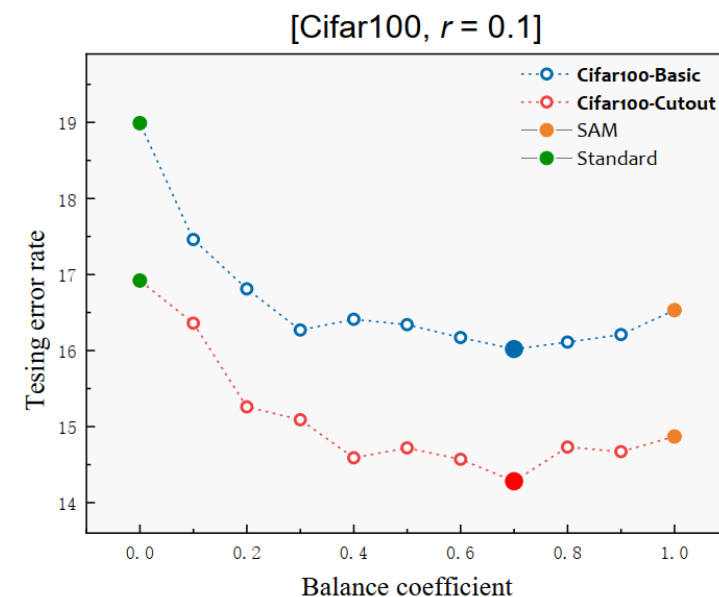
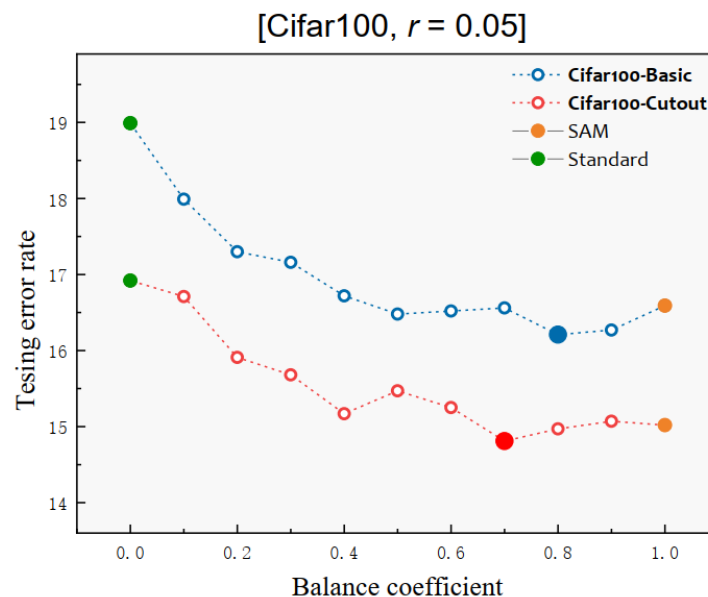
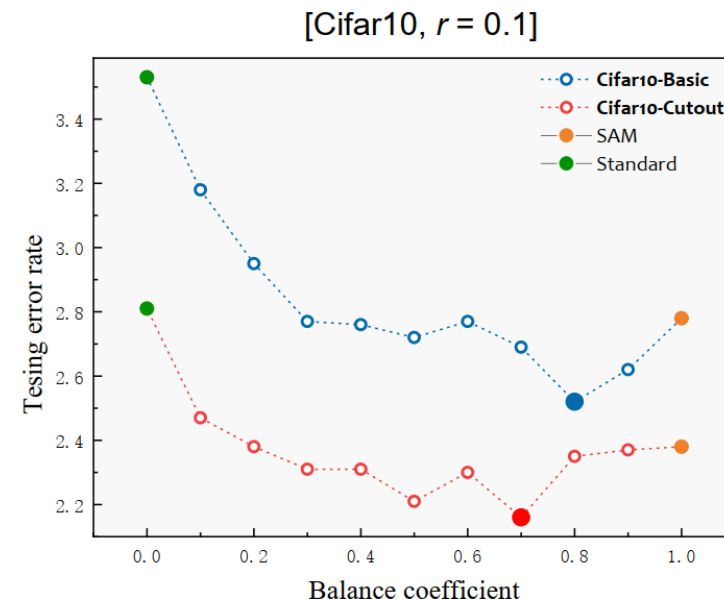
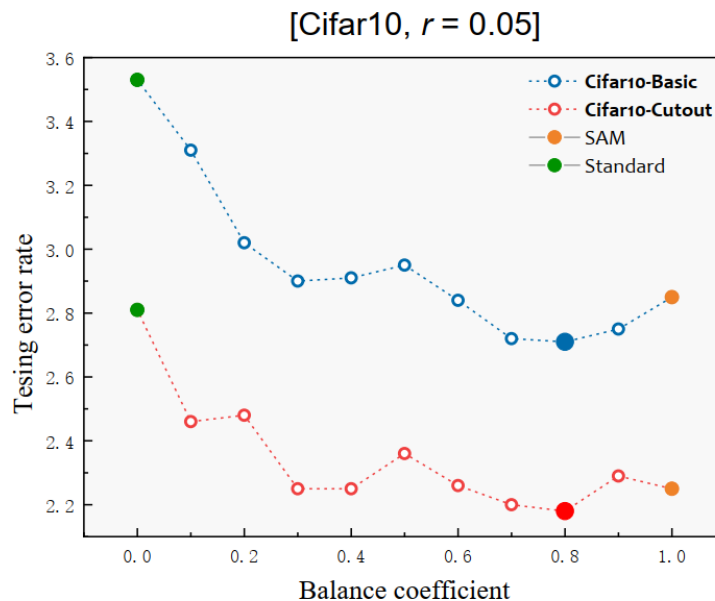
| | Cifar10 | | Cifar100 | |
|----------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| ViT-Ti16 | Basic | Heavy | Basic | Heavy |
| Standard | 15.92 \pm 0.17 | 14.68 \pm 0.14 | 40.21 \pm 0.20 | 38.93 \pm 0.28 |
| SAM | 15.33 \pm 0.18 | 13.77 \pm 0.12 | 38.89 \pm 0.23 | 37.61 \pm 0.19 |
| Ours | 14.75\pm0.17 | 13.52\pm0.21 | 38.58\pm0.27 | 37.15\pm0.21 |
| ViT-S16 | Basic | Heavy | Basic | Heavy |
| Standard | 14.55 \pm 0.14 | 13.31 \pm 0.11 | 38.43 \pm 0.19 | 37.58 \pm 0.22 |
| SAM | 13.91 \pm 0.18 | 12.63 \pm 0.09 | 37.98 \pm 0.23 | 36.77 \pm 0.25 |
| Ours | 13.66\pm0.16 | 12.29\pm0.19 | 37.32\pm0.28 | 36.59\pm0.22 |

- Parameter Study

Figure :

Testing error rate when trained with different hyperparameters. The points colored green are results using the standard scheme, and the points colored orange are results using the SAM scheme.

We recommend to set $\alpha = 0.8$ for achieving the best performance in practice, since most of our training could give the best results at this value.



Thanks for listening!