



# Regularizing Neural Networks via Adversarial Model Perturbation

Yaowei Zheng,<sup>1</sup> Richong Zhang,<sup>1</sup> Yongyi Mao<sup>2</sup>

<sup>1</sup>BDBC and SKLSDE, Beihang University, Beijing, China

<sup>2</sup>School of EECS, University of Ottawa, Ottawa, Canada

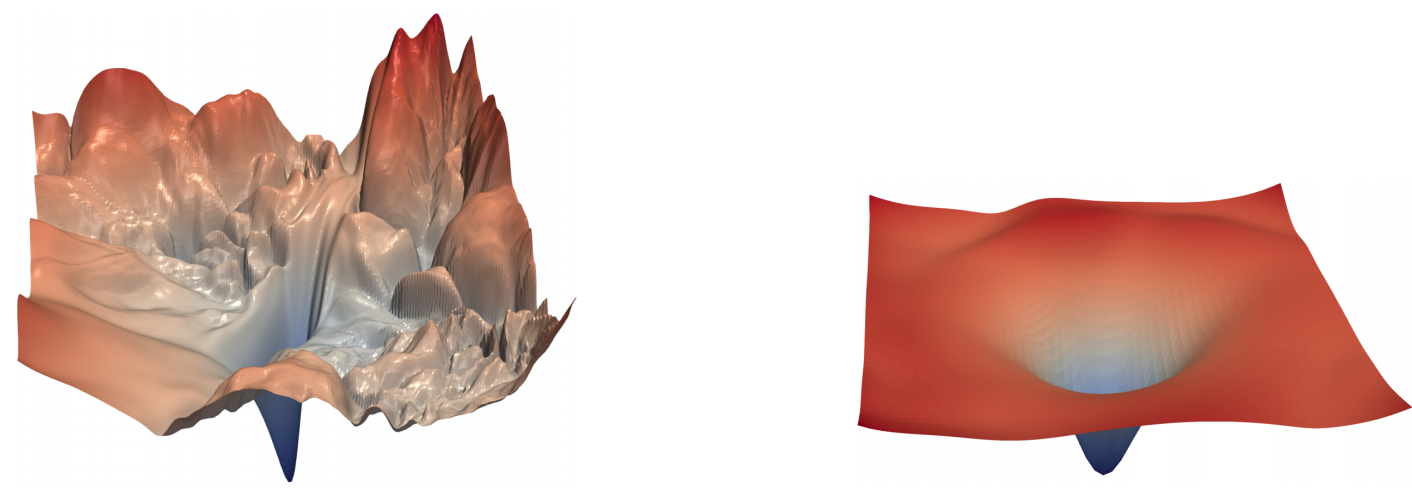


## Summary

**Background:** Effective regularization schemes are highly desired in deep learning for alleviating overfitting and improving generalization.

**Motivation:** Previous work suggested that flat minima can improve generalization both in theoretical and empirical perspectives [1-4].

**Contribution:** We propose Adversarial Model Perturbation (AMP) as a powerful regularization scheme principled by the objective of finding flat minima, which achieves better classification performance and calibration results.



(Better minima are flatter, visualized by Li *et al.*, 2018.)

## Training Algorithm

**Optimization Objective:**

$$\min_{\theta} \max_{\Delta: \|\Delta\| \leq \epsilon} \mathcal{L}_{\text{ERM}}(\theta + \Delta)$$

**Pseudo-code:**

```
while  $\theta$  not converged do
  Initialize perturbation  $\Delta$  with 0;
  for  $n \leftarrow 1$  to  $N$  do
    Update  $\Delta$  to maximize  $\mathcal{L}_{\text{ERM}}(\theta + \Delta)$  via
    gradient ascent with learning rate  $\zeta$ ;
    if  $\|\Delta\|_2 > \epsilon$  then
      Normalize  $\Delta$  to restrict its norm  $\|\Delta\|_2$  to  $\epsilon$ ;
    Update  $\theta$  to minimize  $\mathcal{L}_{\text{ERM}}(\theta + \Delta)$  via gradient
    descent with learning rate  $\eta$ ;
```

**PyTorch Implementation:**

```
from amp import AMP
opt = AMP(params, lr=0.1, epsilon=0.5)
for inputs, targets in dataset:
  def closure():
    opt.zero_grad()
    outputs = model(inputs)
    loss = loss_fn(outputs, targets)
    loss.backward()
  return outputs, loss
outputs, loss = opt.step(closure)
```

**Official Code:**

<https://github.com/hiyouga/AMP-Regularizer>

## Method

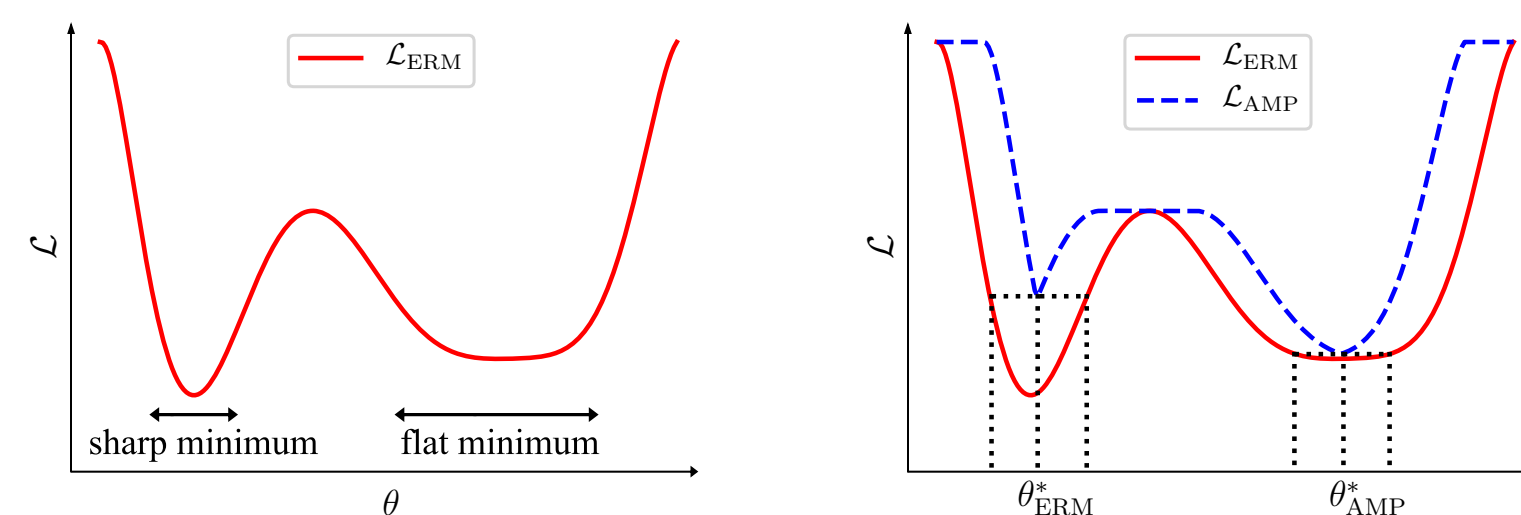
**AMP: Adversarial Model Perturbation**

We derive an AMP loss from the empirical risk (ERM loss). Minimizing such an AMP loss can penalize the sharp minima by adversarially perturbing the model parameters.

$$\mathcal{L}_{\text{ERM}}(\theta) := \frac{1}{|D|} \sum_{(x, y) \in D} \ell(x, y; \theta)$$

$$\mathcal{L}_{\text{AMP}}(\theta) := \max_{\Delta: \|\Delta\| \leq \epsilon} \mathcal{L}_{\text{ERM}}(\theta + \Delta)$$

It can be seen as an analogue of a “max-pooling” operation on the empirical risk on each point in the parameter space.



## Theoretical Justification

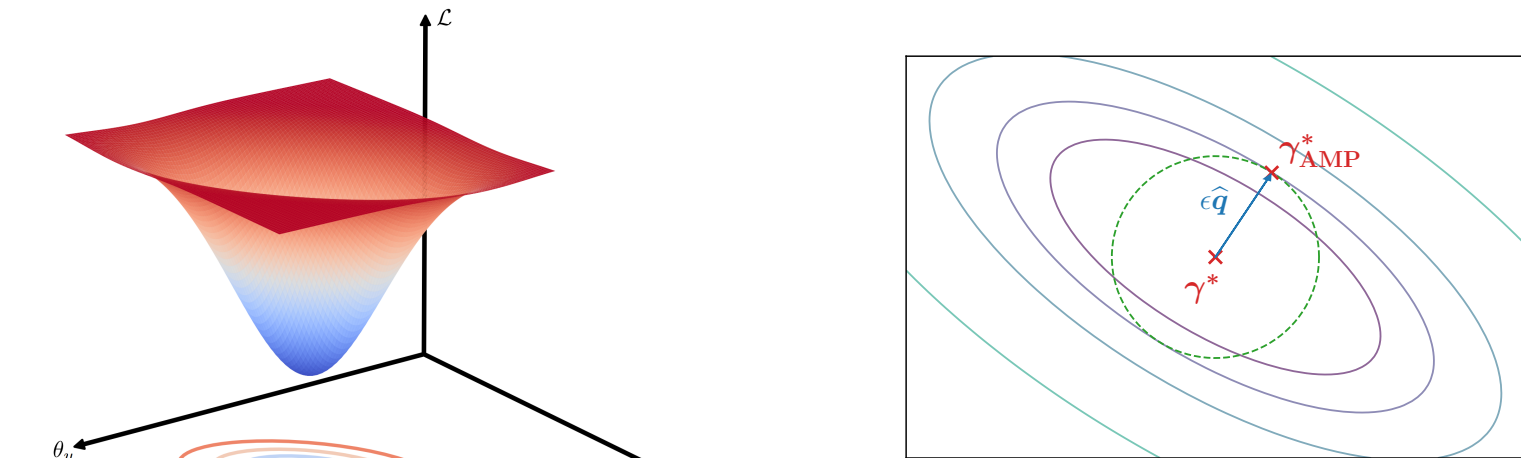
**AMP Finds Flatter Minima:**

We assume that the loss surface of each local minimum in  $\mathcal{L}_{\text{ERM}}$  can be *locally* approximated as an inverted Gaussian surface  $\gamma$  with a mean vector  $\mu$  and a covariance matrix  $\kappa$ . Under the locally Gaussian assumption, the empirical risk  $\gamma(\theta; \mu, \kappa, A, C)$  is minimized when  $\theta = \mu$  and the minimum value is  $\gamma^*(\mu, \kappa, A, C) = C - A$ . The minimum value of the AMP loss is the empirical risk at the location in the narrowest principal direction of the cross-section of the loss surface:

$$\gamma_{\text{AMP}}^*(\mu, \kappa, A, C) = C - A \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)$$

where  $\sigma^2$  is the smallest eigenvalue of  $\kappa$ .

It is clear that  $\gamma_{\text{AMP}}^*$  although related to the minimum value of empirical risk  $\gamma^*$ , it also takes into account the curvature of the surface around the local minimum. Thus we can find a flatter minimum through minimizing  $\mathcal{L}_{\text{AMP}}$ .



**AMP Regularizes Gradient Norm:**

Consider that  $N = 1$ . The AMP training is equivalent to ERM training with an additional term:

$$\tilde{\mathcal{J}}_{\text{ERM}}(\theta) := \mathcal{J}_{\text{ERM}}(\theta) + \Omega(\theta)$$

where

$$\Omega(\theta) := \begin{cases} \zeta \|\nabla_{\theta} \mathcal{J}_{\text{ERM}}(\theta)\|_2^2, & \|\zeta \nabla_{\theta} \mathcal{J}_{\text{ERM}}(\theta)\|_2 \leq \epsilon \\ \epsilon \|\nabla_{\theta} \mathcal{J}_{\text{ERM}}(\theta)\|_2, & \|\zeta \nabla_{\theta} \mathcal{J}_{\text{ERM}}(\theta)\|_2 > \epsilon \end{cases}$$

## Experiments

**Results on Image Classification Benchmarks:**

PreActResNet18	Test Error (%)	Test NLL
ERM	2.95±0.063	0.166±0.004
Dropout	2.80±0.065	0.156±0.012
Label Smoothing	2.78±0.087	0.998±0.002
Flooding	2.84±0.047	0.130±0.003
MixUp	2.74±0.044	0.146±0.004
Adv. Training	2.77±0.080	0.151±0.018
RMP	2.93±0.066	0.161±0.010
AMP	<b>2.30±0.025</b>	<b>0.096±0.002</b>

(a) SVHN

PreActResNet18	Test Error (%)	Test NLL
ERM	5.02±0.212	0.239±0.009
Dropout	4.86±0.148	0.223±0.009
Label Smoothing	4.85±0.115	1.038±0.003
Flooding	4.97±0.082	0.166±0.003
MixUp	4.09±0.117	0.198±0.004
Adv. Training	4.99±0.085	0.247±0.006
RMP	4.97±0.167	0.239±0.008
AMP	<b>3.97±0.091</b>	<b>0.129±0.003</b>

(b) CIFAR-10

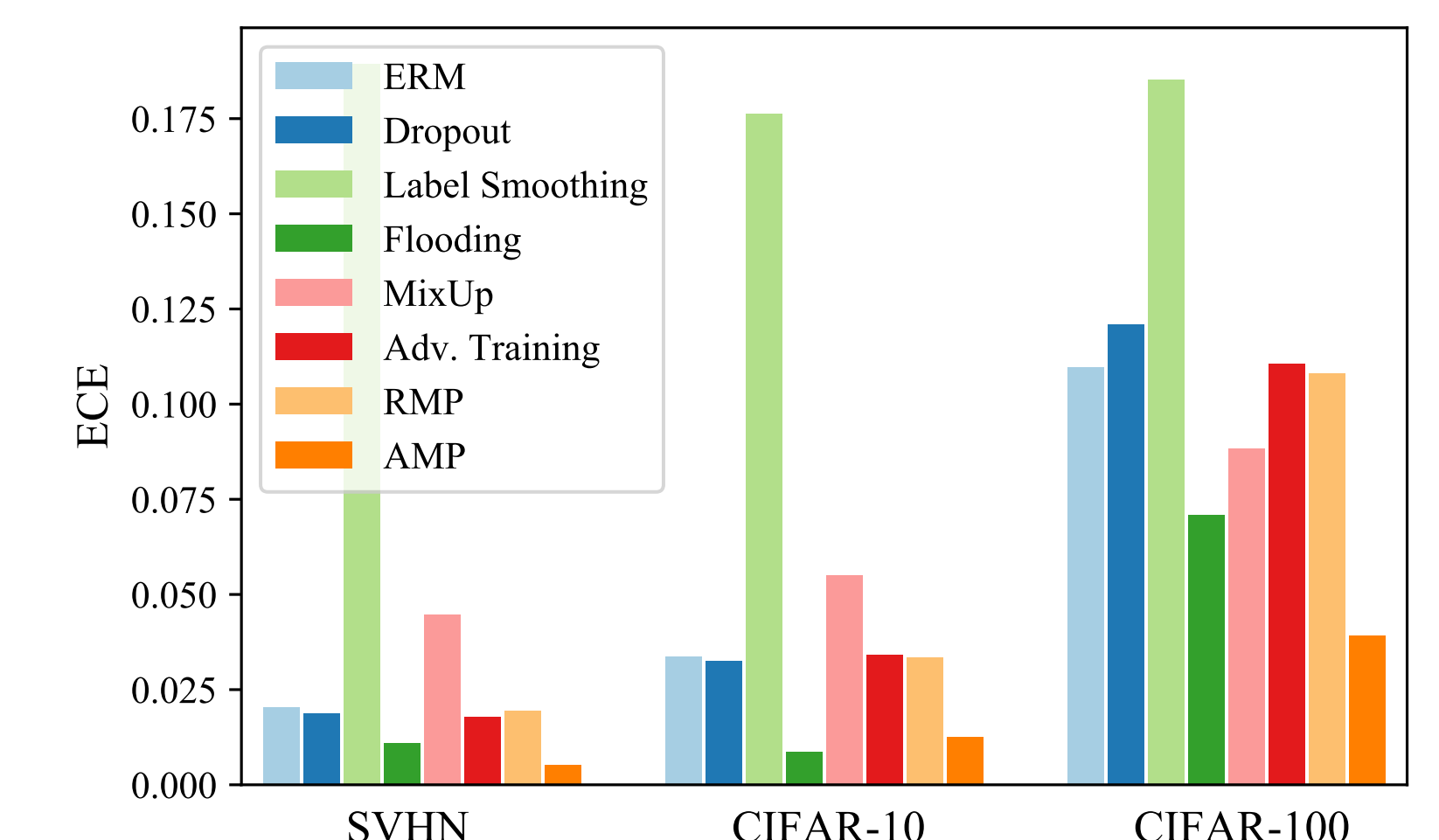
PreActResNet18	Test Error (%)	Test NLL
ERM	24.31±0.303	1.056±0.013
Dropout	24.48±0.351	1.110±0.021
Label Smoothing	22.07±0.256	2.099±0.005
Flooding	24.50±0.234	0.950±0.011
MixUp	21.78±0.210	0.910±0.007
Adv. Training	25.23±0.229	1.110±0.012
RMP	24.28±0.138	1.059±0.011
AMP	<b>21.51±0.308</b>	<b>0.774±0.016</b>

(c) CIFAR-100

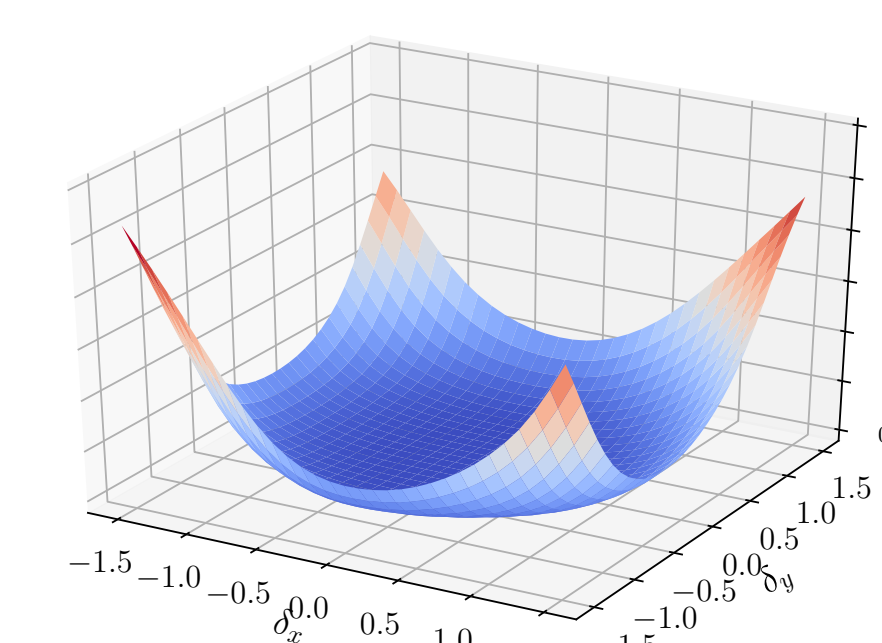
**Improvement over Various Data Augmentation Techniques:**

Test Error (%)	Augment	WideResNet-28-10		PyramidNet-164-270	
Dataset	Technique	ERM	AMP	ERM	AMP
SVHN	Vanilla	2.57±0.067	<b>2.19±0.036</b>	2.47±0.034	<b>2.11±0.041</b>
	Cutout	2.27±0.085	<b>1.83±0.018</b>	2.19±0.021	<b>1.82±0.023</b>
	AutoAug	1.91±0.059	<b>1.61±0.024</b>	1.80±0.044	<b>1.35±0.056</b>
CIFAR-10	Vanilla	3.87±0.167	<b>3.00±0.059</b>	3.60±0.197	<b>2.75±0.040</b>
	Cutout	3.38±0.081	<b>2.67±0.043</b>	2.83±0.102	<b>2.27±0.034</b>
	AutoAug	2.78±0.134	<b>2.32±0.097</b>	2.49±0.128	<b>1.98±0.062</b>
CIFAR-100	Vanilla	19.17±0.270	<b>17.33±0.110</b>	17.13±0.210	<b>15.09±0.092</b>
	Cutout	18.12±0.114	<b>16.04±0.071</b>	16.45±0.136	<b>14.34±0.153</b>
	AutoAug	17.79±0.185	<b>14.95±0.088</b>	15.43±0.269	<b>13.36±0.245</b>

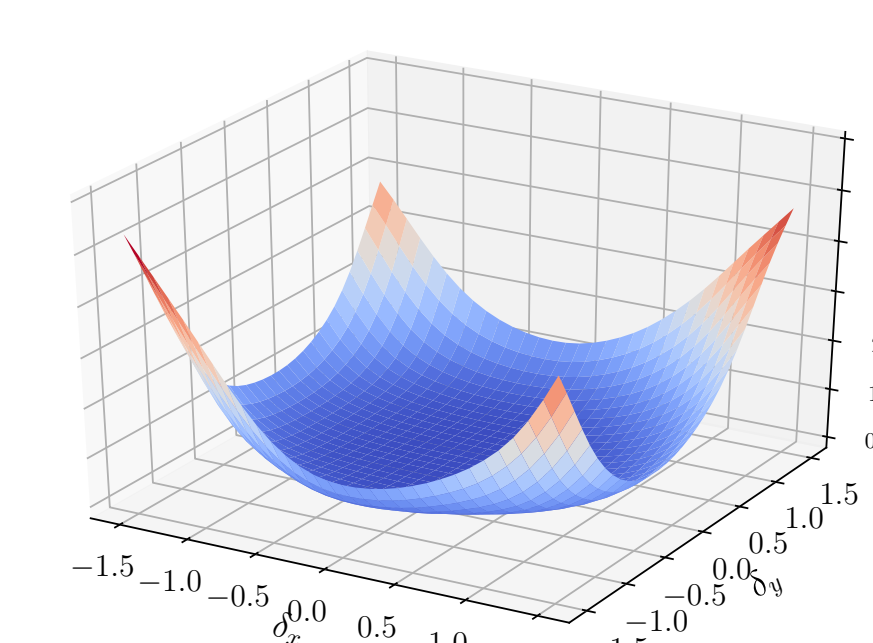
**Calibration Results: (lower is better)**



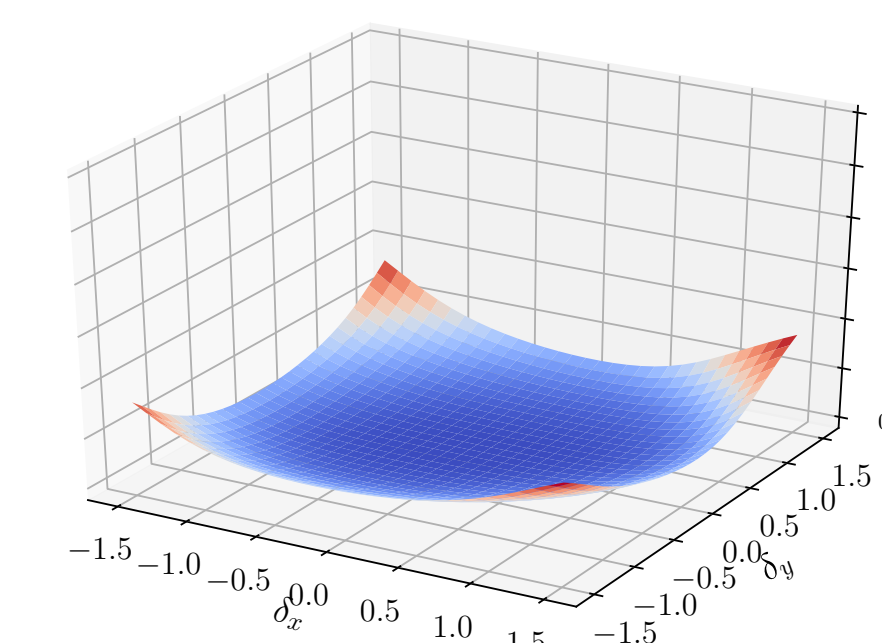
**Flatness of the Selected Minima:**



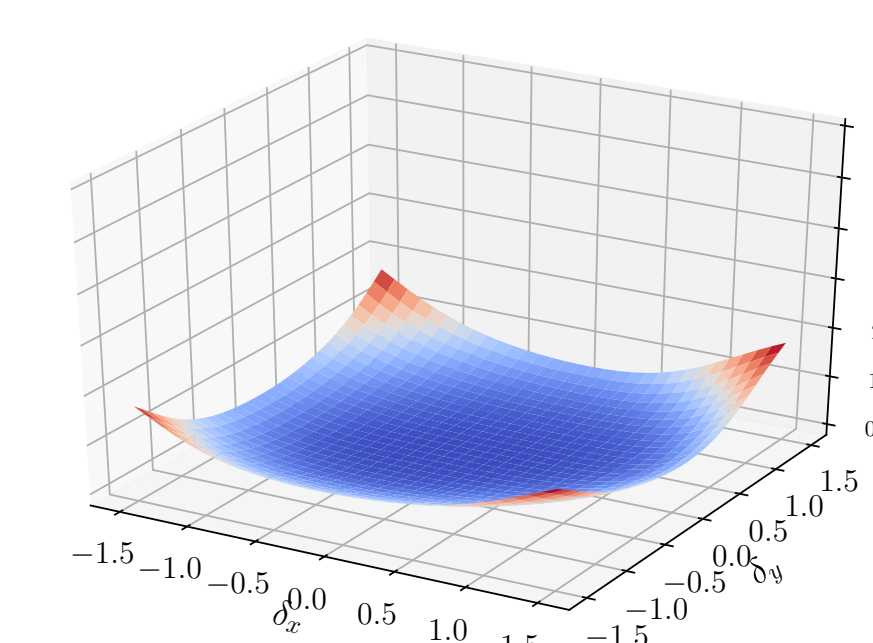
ERM Training Loss



ERM Test Loss

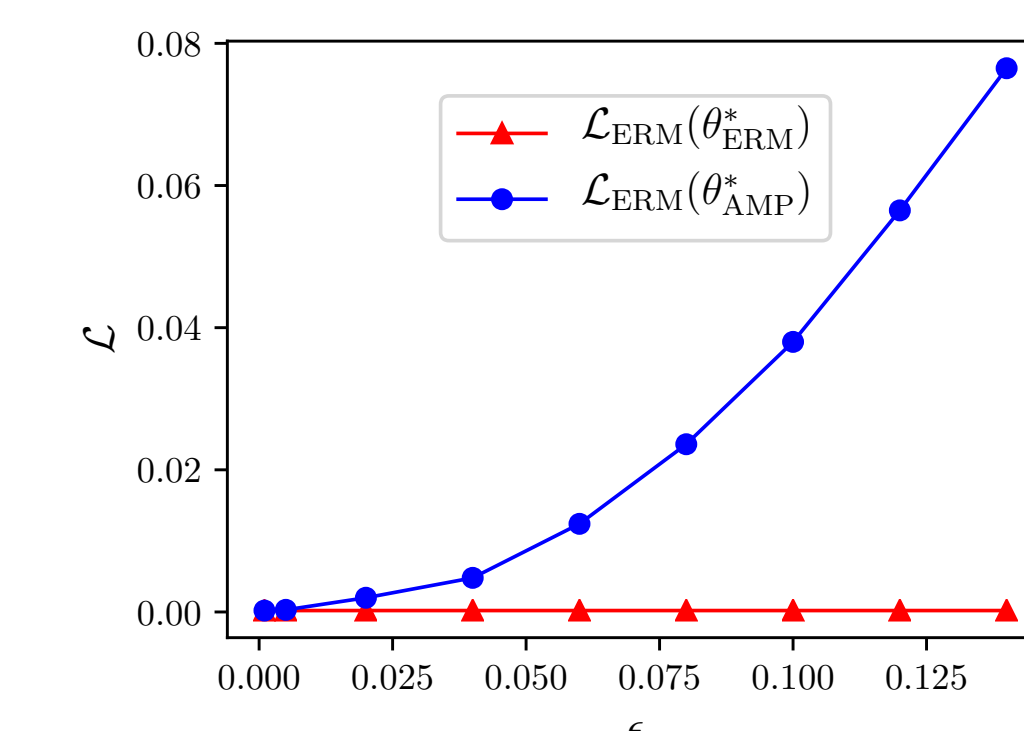


AMP Training Loss

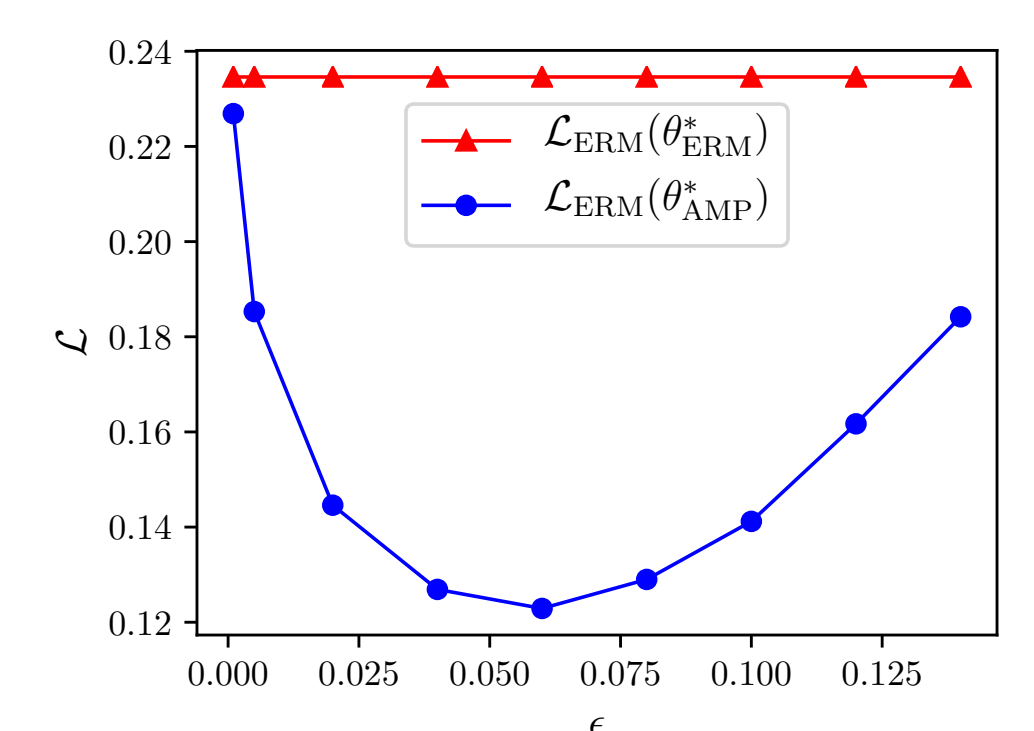


AMP Test Loss

**Loss Values with Varying Perturbation Size:**



(a) CIFAR-10 Training Set



(b) CIFAR-10 Test Set

**References:**

- [1] Hochreiter *et al.* Flat minima. Neural Computation, 9(1):1-42, 1997.
- [2] Keskar *et al.* On large-batch training for deep learning: Generalization gap and sharp minima. In ICLR, 2017.
- [3] Li *et al.* Visualizing the loss landscape of neural nets. In NeurIPS, 2018.
- [4] Foret *et al.* Sharpness-aware minimization for efficiently improving generalization. In ICLR, 2021.