



# Towards Understanding the Importance of Noise in Training Neural Networks

Mo Zhou<sup>2</sup>, Tianyi Liu<sup>1</sup>, Yan Li<sup>1</sup>, Dachao Lin<sup>2</sup>, Enlu Zhou<sup>1</sup>, Tuo Zhao<sup>1</sup>  
<sup>1</sup> Georgia Institute of Technology <sup>2</sup> Peking University



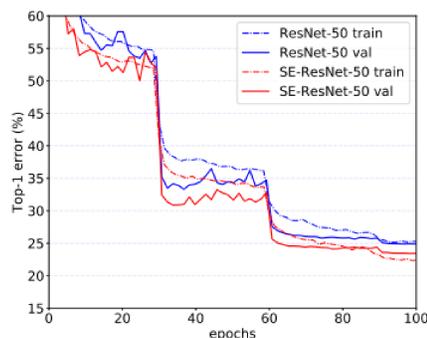
## Background

Challenges of Deep Neural Networks (DNNs):

- Nonconvex: Saddle Points, Spurious Optima;
- Computationally intractable;
- Serious overfitting and curse of dimensionality.

Success of First Order Methods:

- Existing Results:
  - Escape strict saddle and converge to optima:
    - Gradient Descent (GD): Lee et al., 2016, etc.;
    - Stochastic Gradient Descent (SGD): Dauphin et al., 2014, etc.
- Practitioners' Choice: Step Size Annealing (SSA)



- Our Empirical Observations:

	Generalization	Optimal	Noise Level
GD	Bad	Sharp	No
SGD w./ very small Step Size	Bad	Sharp	Very Small
SGD w./ SSA	Good	Flat	Stagewise Decreasing

- Not all local optima generalize well.
- Noise can help select good optima.

**Question: How does noise help train neural networks in the presence of bad optima?**

## Challenges

- Beyond Technical Limit:
  - General NNs: Complex Nonconvex Landscape;
  - SGD: Noise Structure: Distribution and Dependency.
- We Study:
  - Two-Layer Nonoverlapping Convolutional NNs (CNNs):
    1. A non-trivial spurious local optimum;
    2. GD gets trapped with constant probability ( $\frac{1}{4} \sim \frac{3}{4}$ );
  - Perturbed Gradient Descent with Noise Annealing:
    1. Independent injected Uniform noise;
    2. Imitate the behavior of SGD.

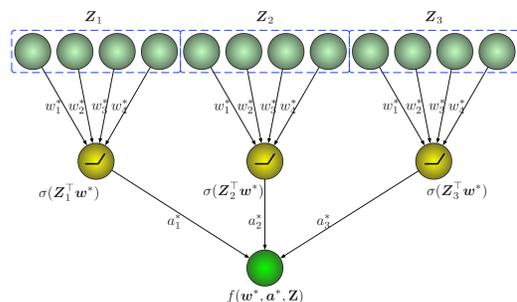
**A non-trivial example provides new insights!**

## Two-layer Nonoverlapping CNNs

- Teacher Network Model:

$$f(w^*, a^*, \mathbf{Z}) = \sum_{j=1}^k a_j^* \sigma(\mathbf{Z}_j^\top w^*),$$

- $\|w^*\|_2 = 1, w \in \mathbb{R}^p, a \in \mathbb{R}^k,$
- $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$  with  $\mathbf{Z}_j$ 's i.i.d.  $N(\mathbf{0}, \mathbf{I}),$
- $\sigma(\cdot) = \max\{\cdot, 0\}.$



- Nonconvex Optimization:

$$(\hat{w}, \hat{a}) = \operatorname{argmin}_{w, a} \mathcal{L}(w, a) \quad \text{subject to} \quad \|w\|_2 = 1,$$

where  $\mathcal{L}(w, a) = \mathbb{E}_{\mathbf{Z}}(f(w^*, a^*, \mathbf{Z}) - f(w, a, \mathbf{Z}))^2.$

- A nontrivial spurious local optimum exists! (Du et al., 2017)

$$\bar{w} = -w^*, \quad \bar{a} = (\mathbf{1}\mathbf{1}^\top + (\pi - 1)\mathbf{I})^{-1}(\mathbf{1}\mathbf{1}^\top - \mathbf{I})a^*.$$

## Perturbed Gradient Descent (P-GD)

- Initialization:  $a_0 \in \mathbb{B}_0\left(\frac{\|\mathbf{1}^\top a^*\|}{\sqrt{k}}\right)$  and  $w_0 \in \mathbb{S}_0(1).$

- At the  $t$ -th iteration, we independently sample

$$\epsilon_t \sim \text{Unif}(\mathbb{B}^p(\rho_w)), \quad \xi_t \sim \text{Unif}(\mathbb{B}^k(\rho_a)),$$

and further take

$$\tilde{w}_t = w_t + \xi_t, \quad \tilde{a}_t = a_t + \epsilon_t.$$

- We then update  $w$  and  $a$  by

$$a_{t+1} = a_t - \eta \nabla_a \mathcal{L}_t(\tilde{w}_t, \tilde{a}_t, \mathbf{Z}^{(t)}),$$

$$w_{t+1} = \Pi_{\mathbb{S}(1)}(w_t - \eta(\mathbf{I} - w_t w_t^\top) \nabla_w \mathcal{L}_t(\tilde{w}_t, \tilde{a}_t, \mathbf{Z}^{(t)})),$$

where  $\mathcal{L}(w, a, \mathbf{Z}) = \mathbb{E}_{\mathbf{Z}}(f(w^*, a^*, \mathbf{Z}) - f(w, a, \mathbf{Z}))^2.$

### Noise Annealing:

Noise level schedule:  $\{\rho_w^{(s)}\}_{s=1}^S$  and  $\{\rho_a^{(s)}\}_{s=1}^S$

- Multi-Epoch:

$$s\text{-th initialization} \leftarrow (s-1)\text{-th output}.$$

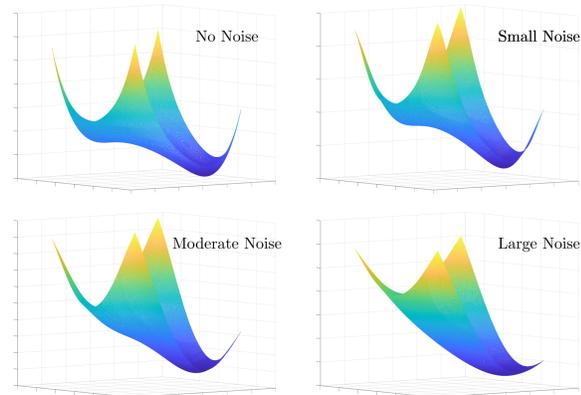
- Noise Annealing:

$$\rho_w^{(s)} < \rho_w^{(s-1)} \quad \text{and} \quad \rho_a^{(s)} < \rho_a^{(s-1)}$$

## Convolutional Effects

- P-GD essentially solves:

$$\min_{w, a} \mathbb{E}_{\epsilon, \xi} \mathcal{L}(w + \epsilon, a + \xi) \quad \text{subject to} \quad \|w\|_2 = 1.$$



- Partial Dissipativity Condition: Let  $\tilde{w} = w + \epsilon$  and  $\tilde{a} = a + \xi.$  For  $(w, a) \in \mathcal{U} \subseteq \mathbb{S}^p(1) \times \mathbb{B}^k(R),$  we have:

$$C1: \langle -\mathbb{E}_{\xi, \epsilon} (\mathbf{I} - w w^\top) \nabla_w \mathcal{L}(\tilde{w}, \tilde{a}), w^* - w \rangle \geq c_w \|w - w^*\|_2^2 - \gamma_w,$$

$$C2: \langle -\mathbb{E}_{\xi, \epsilon} \nabla_a \mathcal{L}(\tilde{w}, \tilde{a}), a^* - a \rangle \geq c_a \|a - a^*\|_2^2 - \gamma_a.$$

- Technical Challenges:

- C1 and C2 do NOT globally hold;
- C1 and C2 do NOT necessarily hold at the same time;
- C1 and C2 vary as the noise levels vary.

## Epoch I: Escaping Spurious Local Optima

Partial Dissipativity Condition: With large noise,

- C2 holds and C1 does not hold around the initialization.  $\Rightarrow$  P-GD reduces the optimization error of  $a.$
- Reducing error of  $a. \Rightarrow$  C1 holds.  $\Rightarrow$  P-GD improves  $w.$
- The output solution is far away from  $(w^*, a^*).$

**Theorem 1.** Suppose  $\rho_w^0 = C_w^0 k p^2 \geq 1$  and  $\rho_a^0 = C_a^0.$  For any  $\delta \in (0, 1),$  we choose step size

$$\eta = O\left(\left(k^4 p^6 \cdot \max\left\{1, p \log \frac{1}{\delta}\right\}\right)^{-1}\right).$$

Then with probability at least  $1 - \delta,$  we have

$$0 < m_a \leq a_t^\top a^* \leq M_a \quad \text{and} \quad \angle(w_t, w^*) \leq \frac{5}{12} \pi$$

for all  $T_1 \leq t \leq O(\eta^{-2}),$  where  $m_a, M_a$  are some constants, and

$$T_1 = O\left(p k / \eta \log(1/\eta) \log(1/\delta) \|a^*\|_2^2\right).$$

## Epoch II: Converging to Global Optima

Partial Dissipativity Condition: With small noise,

- C1 and C2 jointly hold;
- $\gamma_w = 0$  and  $\gamma_a$  decreases.

**Theorem 2.** For any  $\gamma > 0,$  we choose  $\rho_w^1 \leq C_w^1 \frac{\gamma}{k p} < 1$  and  $\rho_a^1 \leq C_a^1$  for some constants  $C_w^1$  and  $C_a^1.$  For any  $\delta \in (0, 1),$  we choose step size

$$\eta = O\left(\left(\max\left\{k^4 p^6, \frac{k^2 p}{\gamma}\right\} \max\left\{1, p \log \frac{1}{\gamma} \log \frac{1}{\delta}\right\}\right)^{-1}\right).$$

Then with probability at least  $1 - \delta,$  we have

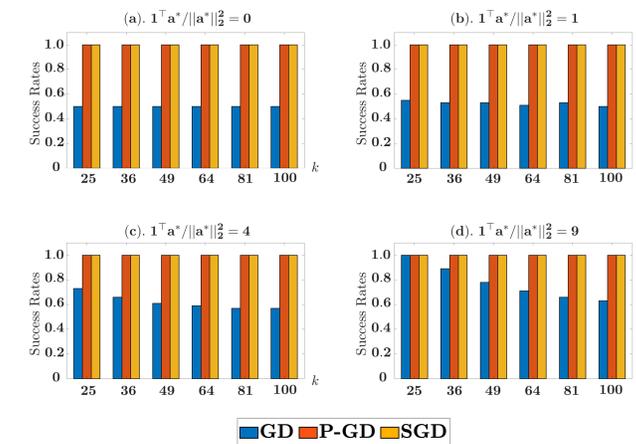
$$\|w_t - w^*\|_2^2 \leq \gamma \quad \text{and} \quad \|a_t - a^*\|_2^2 \leq \gamma$$

for any  $t$ 's such that  $T_2 \leq t \leq T = O(\eta^{-2}),$  where

$$T_2 = O\left(p / \eta \log 1/\gamma \log 1/\delta \|a^*\|_2^2\right).$$

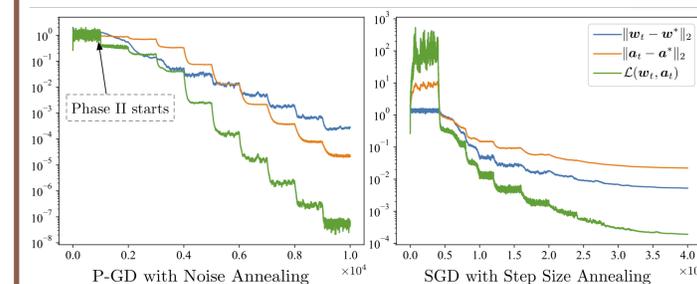
## Experiments

- Success Rates with  $p = 6$  and Varing  $k$  and  $a^*:$



- P-GD, SGD 100% converge to the global optimum.
- GD can get trapped with probability 50%.

- Empirical Convergence:



- SGD shows similar patterns to P-GD.