

Towards Understanding the Importance of Noise in Training Neural Networks

Mo Zhou[#], **Tianyi Liu**[†], Yan Li[†], Dachao Lin[#],
Enlu Zhou[†] and Tuo Zhao[†]

[†]Georgia Tech and [#]Peking University

June. 12, 2019

A Natural Question:

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

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challenges

General Neural Networks (NNs)

- Complex nonconvex landscape;
- Beyond our technical limit.

We Study: Two-Layer Nonoverlapping Convolutional NNs:

- Non-trivial spurious local optimum (does not generalize);
- GD with random initialization gets trapped with constant probability (at least $\frac{1}{4}$, can be $\frac{3}{4}$ in the worst case);
- Simple structure that is technically manageable.

Challengess

Stochastic Gradient Descent

- Complex distribution of noise;
- Dependency on iterates.

We Study: Perturbed Gradient Descent with Noise Annealing:

- Independent injected noise;
- Uniform distribution;
- Imitate the behavior of SGD.

A non-trivial example provides new insights!

Challengess

Stochastic Gradient Descent

- Complex distribution of noise;
- Dependency on iterates.

We Study: Perturbed Gradient Descent with Noise Annealing:

- Independent injected noise;
- Uniform distribution;
- Imitate the behavior of SGD.

A non-trivial example provides new insights!

Challengess

Stochastic Gradient Descent

- Complex distribution of noise;
- Dependency on iterates.

We Study: Perturbed Gradient Descent with Noise Annealing:

- Independent injected noise;
- Uniform distribution;
- Imitate the behavior of SGD.

A non-trivial example provides new insights!

Challengess

Stochastic Gradient Descent

- Complex distribution of noise;
- Dependency on iterates.

We Study: Perturbed Gradient Descent with Noise Annealing:

- Independent injected noise;
- Uniform distribution;
- Imitate the behavior of SGD.

A non-trivial example provides new insights!

Challengess

Stochastic Gradient Descent

- Complex distribution of noise;
- Dependency on iterates.

We Study: Perturbed Gradient Descent with Noise Annealing:

- Independent injected noise;
- Uniform distribution;
- Imitate the behavior of SGD.

A non-trivial example provides new insights!

Two-layer Nonoverlapping CNNs

- Teacher Network Model:

$$f(\mathbf{w}^*, \mathbf{a}^*, \mathbf{Z}) = \sum_{j=1}^k a_j^* \sigma(\mathbf{Z}_j^\top \mathbf{w}^*) \quad \text{with} \quad \|\mathbf{w}^*\|_2 = 1,$$

where $\mathbf{w} \in \mathbb{R}^p$, $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$ with \mathbf{Z}_j 's are independently sampled from $N(\mathbf{0}, \mathbf{I})$, and $\sigma(\cdot) = \max\{\cdot, 0\}$.

- Nonconvex Optimization:

$$(\hat{\mathbf{w}}, \hat{\mathbf{a}}) = \arg \min_{\mathbf{w}, \mathbf{a}} \mathcal{L}(\mathbf{w}, \mathbf{a}) \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1,$$

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

- A nontrivial spurious local optimum $(\bar{\mathbf{w}}, \bar{\mathbf{a}})$ exists!

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

Two-layer Nonoverlapping CNNs

- Teacher Network Model:

$$f(\mathbf{w}^*, \mathbf{a}^*, \mathbf{Z}) = \sum_{j=1}^k a_j^* \sigma(\mathbf{Z}_j^\top \mathbf{w}^*) \quad \text{with} \quad \|\mathbf{w}^*\|_2 = 1,$$

where $\mathbf{w} \in \mathbb{R}^p$, $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$ with \mathbf{Z}_j 's are independently sampled from $N(\mathbf{0}, \mathbf{I})$, and $\sigma(\cdot) = \max\{\cdot, 0\}$.

- Nonconvex Optimization:

$$(\hat{\mathbf{w}}, \hat{\mathbf{a}}) = \arg \min_{\mathbf{w}, \mathbf{a}} \mathcal{L}(\mathbf{w}, \mathbf{a}) \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1,$$

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

- A nontrivial spurious local optimum $(\bar{\mathbf{w}}, \bar{\mathbf{a}})$ exists!

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

Two-layer Nonoverlapping CNNs

- Teacher Network Model:

$$f(\mathbf{w}^*, \mathbf{a}^*, \mathbf{Z}) = \sum_{j=1}^k a_j^* \sigma(\mathbf{Z}_j^\top \mathbf{w}^*) \quad \text{with} \quad \|\mathbf{w}^*\|_2 = 1,$$

where $\mathbf{w} \in \mathbb{R}^p$, $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$ with \mathbf{Z}_j 's are independently sampled from $N(\mathbf{0}, \mathbf{I})$, and $\sigma(\cdot) = \max\{\cdot, 0\}$.

- Nonconvex Optimization:

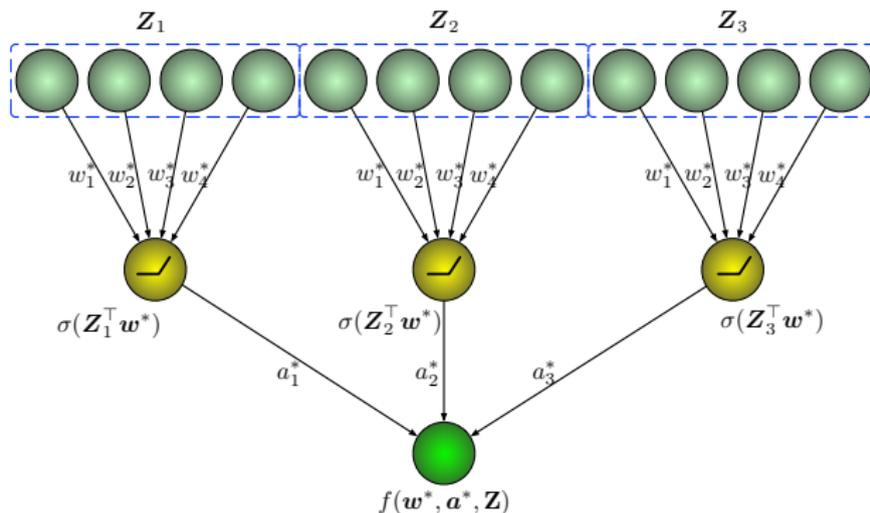
$$(\hat{\mathbf{w}}, \hat{\mathbf{a}}) = \arg \min_{\mathbf{w}, \mathbf{a}} \mathcal{L}(\mathbf{w}, \mathbf{a}) \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1,$$

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

- A nontrivial spurious local optimum $(\bar{\mathbf{w}}, \bar{\mathbf{a}})$ exists!

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

Teacher Network Model



Perturbed Gradient Descent (P-GD)

Initialization: $a_0 \in \mathbb{B}_0 \left(|\mathbf{1}^\top a^*| / \sqrt{k} \right)$ 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 $\mathbf{Z}^{(t)} = [\mathbf{Z}_1^{(t)}, \dots, \mathbf{Z}_t^{(t)}]$ with $\mathbf{Z}_j^{(t)} \sim N(0, \mathbf{I})$, 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 $\ell(w, a, \mathbf{Z}) = \frac{1}{2}(f(w^*, a^*, \mathbf{Z}) - f(w, a, \mathbf{Z}))^2$.

Perturbed Gradient Descent (P-GD)

Initialization: $a_0 \in \mathbb{B}_0 \left(|\mathbf{1}^\top a^*| / \sqrt{k} \right)$ 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 $\mathbf{Z}^{(t)} = [\mathbf{Z}_1^{(t)}, \dots, \mathbf{Z}_t^{(t)}]$ with $\mathbf{Z}_j^{(t)} \sim N(0, \mathbf{I})$, 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 $\ell(w, a, \mathbf{Z}) = \frac{1}{2} (f(w^*, a^*, \mathbf{Z}) - f(w, a, \mathbf{Z}))^2$.

Perturbed Gradient Descent (P-GD)

Initialization: $a_0 \in \mathbb{B}_0 \left(|\mathbf{1}^\top a^*| / \sqrt{k} \right)$ 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 $\mathbf{Z}^{(t)} = [\mathbf{Z}_1^{(t)}, \dots, \mathbf{Z}_t^{(t)}]$ with $\mathbf{Z}_j^{(t)} \sim N(0, \mathbf{I})$, and further take

$$\tilde{\mathbf{w}}_t = \mathbf{w}_t + \xi_t, \quad \tilde{\mathbf{a}}_t = \mathbf{a}^{(t)} + \epsilon_t.$$

We then update \mathbf{w} and \mathbf{a} by

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

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

where $\ell(\mathbf{w}, \mathbf{a}, \mathbf{Z}) = \frac{1}{2} (f(\mathbf{w}^*, \mathbf{a}^*, \mathbf{Z}) - f(\mathbf{w}, \mathbf{a}, \mathbf{Z}))^2$.

Noise Annealing

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

- Multi-Epoch: At the s -th epoch, we initialize using the output solution of the $(s - 1)$ -th epoch.
- Then we apply P-GD with

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

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

Noise Annealing

Noise level schedule: $\{\rho_{\mathbf{w}}^{(s)}\}_{s=1}^S$ and $\{\rho_{\mathbf{a}}^{(s)}\}_{s=1}^S$

- Multi-Epoch: At the s -th epoch, we initialize using the output solution of the $(s - 1)$ -th epoch.
- Then we apply P-GD with

$$\epsilon_t \sim \text{Unif}(\mathbb{B}^p(\rho_{\mathbf{w}}^{(s)})), \quad \xi_t \sim \text{Unif}(\mathbb{B}^k(\rho_{\mathbf{a}}^{(s)})),$$

where $\rho_{\mathbf{w}}^{(s)} < \rho_{\mathbf{w}}^{(s-1)}$ and $\rho_{\mathbf{a}}^{(s)} < \rho_{\mathbf{a}}^{(s-1)}$

Algorithmic Behaviors

Nonasymptotic Convergence Analysis

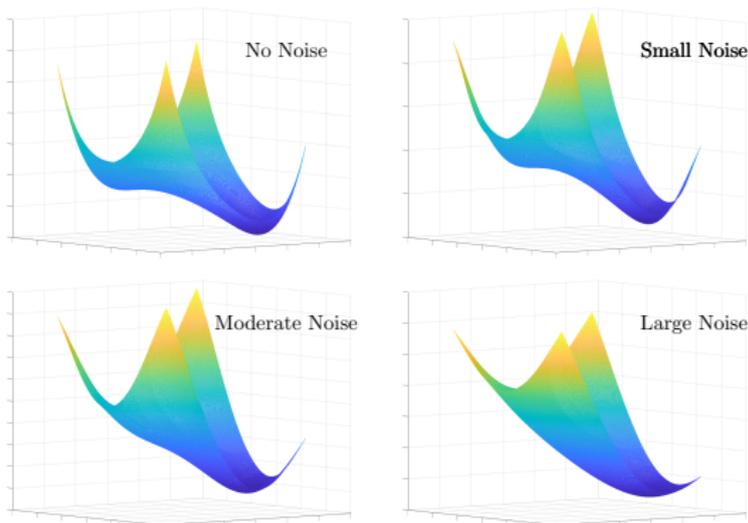
Theorem (Informal)

The theory considers two epochs of P-GD:

- *Epoch I. With **large noise** and a properly chosen step size, P-GD **escapes the local optimum**, while approaching the basin of attraction of the global optimum in polynomial time with high probability;*
- *Epoch II. With **small noise** and a properly chosen step size, P-GD **converges to the global optimum** in polynomial time with high probability.*

“With High Probability”: $\mathbb{P}(A) \geq 1 - \mathcal{O}(\exp(-1/\eta))$.

Convolutional Effects



Remark: P-GD essentially solves:

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

Partially Dissipative Conditions

For $(\mathbf{w}, \mathbf{a}) \in \mathcal{U} \subseteq \mathbb{S}_p(1) \times \mathbb{B}_k(R)$, we have:

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

$$\text{C2: } \langle -\mathbb{E}_{\xi, \epsilon} \nabla_{\mathbf{a}} \mathcal{L}(\tilde{\mathbf{w}}, \tilde{\mathbf{a}}), \mathbf{a}^* - \mathbf{a} \rangle \geq c_{\mathbf{a}} \|\mathbf{a} - \mathbf{a}^*\|_2^2 - \gamma_{\mathbf{a}},$$

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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.

Partially Dissipative Conditions

For $(\mathbf{w}, \mathbf{a}) \in \mathcal{U} \subseteq \mathbb{S}_p(1) \times \mathbb{B}_k(R)$, we have:

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

$$\text{C2: } \langle -\mathbb{E}_{\xi, \epsilon} \nabla_{\mathbf{a}} \mathcal{L}(\tilde{\mathbf{w}}, \tilde{\mathbf{a}}), \mathbf{a}^* - \mathbf{a} \rangle \geq c_{\mathbf{a}} \|\mathbf{a} - \mathbf{a}^*\|_2^2 - \gamma_{\mathbf{a}},$$

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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.

Partially Dissipative Conditions

For $(\mathbf{w}, \mathbf{a}) \in \mathcal{U} \subseteq \mathbb{S}_p(1) \times \mathbb{B}_k(R)$, we have:

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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.

Partially Dissipative Conditions

For $(\mathbf{w}, \mathbf{a}) \in \mathcal{U} \subseteq \mathbb{S}_p(1) \times \mathbb{B}_k(R)$, we have:

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

$$\text{C2: } \langle -\mathbb{E}_{\xi, \epsilon} \nabla_{\mathbf{a}} \mathcal{L}(\tilde{\mathbf{w}}, \tilde{\mathbf{a}}), \mathbf{a}^* - \mathbf{a} \rangle \geq c_{\mathbf{a}} \|\mathbf{a} - \mathbf{a}^*\|_2^2 - \gamma_{\mathbf{a}},$$

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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 the Spurious Local Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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^*) .

Epoch I: Escaping the Spurious Local Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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^*) .

Epoch I: Escaping the Spurious Local Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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^*) .

Epoch I: Escaping the Spurious Local Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

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^*) .

Epoch I: Escaping the Spurious Local Optimum

Theorem

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

$$0 < m_a \leq a_t^\top a^* \leq M_a \quad \text{and} \quad \angle(w_t, w^*) \leq \frac{5}{12}\pi \quad (1)$$

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

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

Remark: (1) is in the basin of attraction of the global optimum.

Epoch I: Escaping the Spurious Local Optimum

Theorem

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

$$0 < m_a \leq a_t^\top a^* \leq M_a \quad \text{and} \quad \angle(w_t, w^*) \leq \frac{5}{12}\pi \quad (1)$$

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

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

Remark: (1) is in the basin of attraction of the global optimum.

Epoch II: Converging to the Global Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

With **small** noise,

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

Epoch II: Converging to the Global Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

With **small** noise,

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

Epoch II: Converging to the Global Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

With **small** noise,

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

Epoch II: Converging to the Global Optimum

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

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

where $\tilde{\mathbf{w}} = \mathbf{w} + \epsilon$ and $\tilde{\mathbf{a}} = \mathbf{a} + \xi$.

With **small** noise,

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

Epoch II: Converging to the Global Optimum

Theorem

For any $\gamma > 0$, we choose $\rho_w^1 \leq C_w^1 \frac{\gamma}{kp} < 1$ and $\rho_a \leq M_a$ for some constant C_w^1 . For any $\delta \in (0, 1)$, we choose step size

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

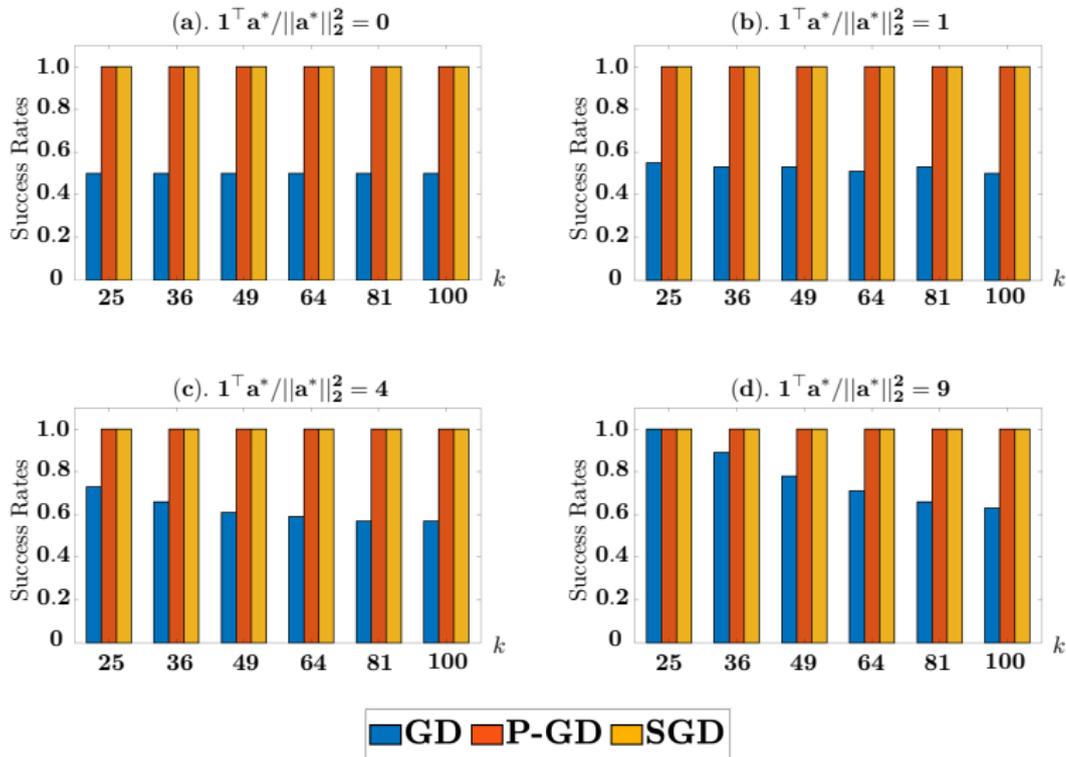
Then with probability at least $1 - \delta$, we have

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

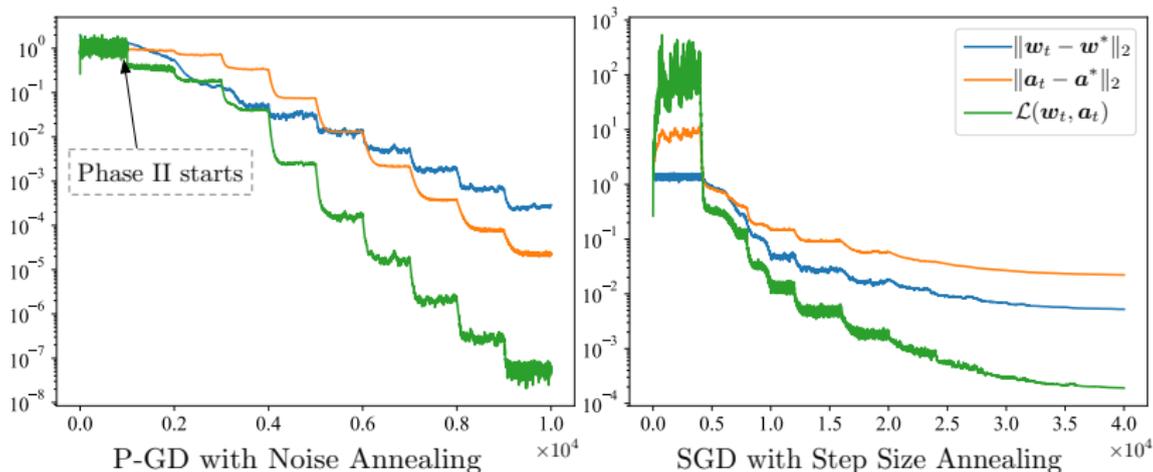
for any t 's such that $T_2 \leq t \leq T = O(\eta^{-2})$, where

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

Experiments: Success Rates ($p=6$)



Experiments: Empirical Convergence



Discussions

- The noise helps escape the local optimum.
- The step size annealing in practice has a similar effect on controlling the noise level.
- P-GD behaves differently for training convolutional weight w and output weight a in the early stage.
- For general deep neural networks, there exist many bad global optima, which cannot generalize. Does SGD escape from them for the same reason?
- To the best of our knowledge, this is the first theoretical result towards justifying the effect of noise in training NNs by SGD-type algorithms in the presence of the spurious optima.

Discussions

- The noise helps escape the local optimum.
- The step size annealing in practice has a similar effect on controlling the noise level.
- P-GD behaves differently for training convolutional weight w and output weight a in the early stage.
- For general deep neural networks, there exist many bad global optima, which cannot generalize. Does SGD escape from them for the same reason?
- To the best of our knowledge, this is the first theoretical result towards justifying the effect of noise in training NNs by SGD-type algorithms in the presence of the spurious optima.

Discussions

- The noise helps escape the local optimum.
- The step size annealing in practice has a similar effect on controlling the noise level.
- P-GD behaves differently for training convolutional weight w and output weight a in the early stage.
- For general deep neural networks, there exist many bad global optima, which cannot generalize. Does SGD escape from them for the same reason?
- To the best of our knowledge, this is the first theoretical result towards justifying the effect of noise in training NNs by SGD-type algorithms in the presence of the spurious optima.

Discussions

- The noise helps escape the local optimum.
- The step size annealing in practice has a similar effect on controlling the noise level.
- P-GD behaves differently for training convolutional weight w and output weight a in the early stage.
- For general deep neural networks, there exist many bad global optima, which cannot generalize. Does SGD escape from them for the same reason?
- To the best of our knowledge, this is the first theoretical result towards justifying the effect of noise in training NNs by SGD-type algorithms in the presence of the spurious optima.

Discussions

- The noise helps escape the local optimum.
- The step size annealing in practice has a similar effect on controlling the noise level.
- P-GD behaves differently for training convolutional weight w and output weight a in the early stage.
- For general deep neural networks, there exist many bad global optima, which cannot generalize. Does SGD escape from them for the same reason?
- To the best of our knowledge, this is the first theoretical result towards justifying the effect of noise in training NNs by SGD-type algorithms in the presence of the spurious optima.

Our Paper:



Poster: Jun. 12 Wed 6:30-9:00 PM Pacific Ballroom No.26

Thank You! Questions?