

# On the Noisy Gradient Descent that Generalizes as SGD

**Jingfeng Wu**, Wenqing Hu, Haoyi Xiong, Jun Huan,  
Vladimir Braverman, Zhanxing Zhu


Johns Hopkins University, Missouri University of Science and Technology,  
Baidu Research, Styling AI, Peking University

# Stochastic gradient descent (SGD)

Loss function

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(x_i; \theta)$$

SGD

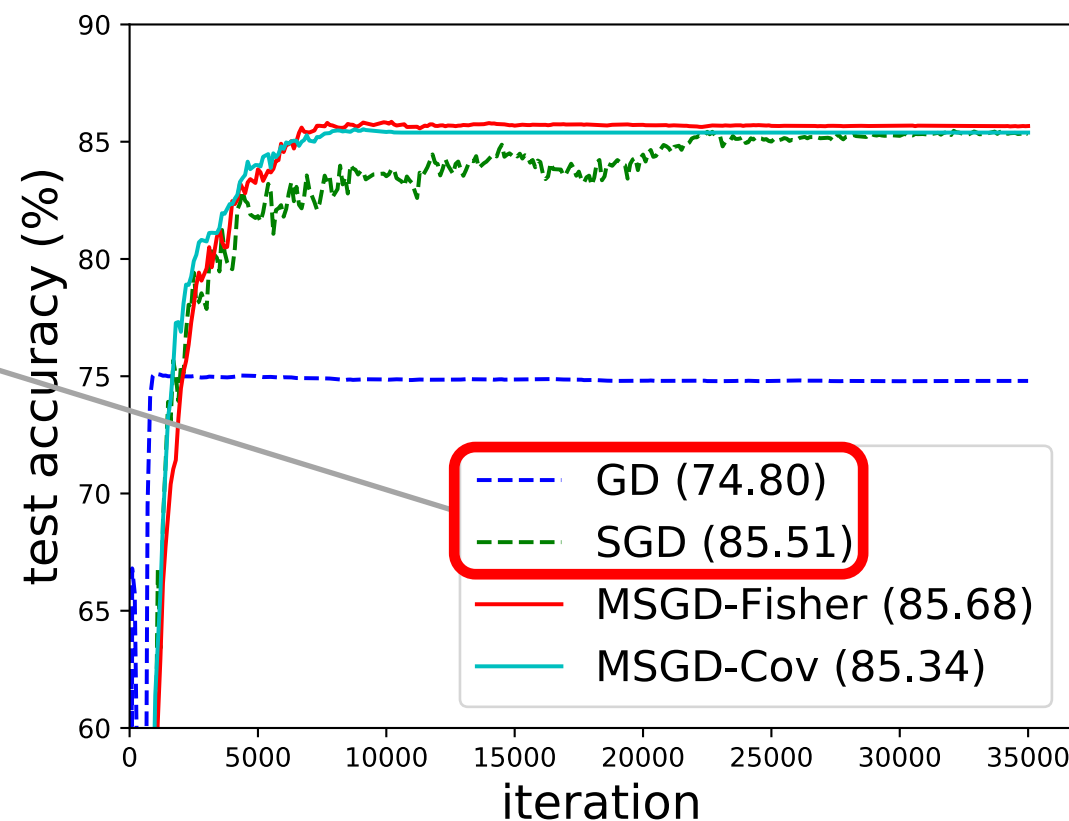

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \underbrace{\tilde{g}(\theta_t)}_{\substack{\text{GD} \\ \frac{1}{b} \sum_{i \in B_t} \nabla_{\theta} \ell(x_i; \theta_t)}} \\ &= \underbrace{\theta_t - \eta \nabla_{\theta} L(\theta_t)}_{\text{GD}} - \underbrace{\eta (\tilde{g}(\theta_t) - \nabla_{\theta} L(\theta_t))}_{v_{\text{sgd}}(\theta_t)} \end{aligned}$$

(unbiased) gradient noise

# Noise matters!

SGD >> GD

- How? <= Still open...
- Which? <= This work!



CIFAR-10, ResNet-18, w/o weight decay,  
w/o data augmentation

# Which noise matters?

$$v_{\text{sgd}}(\theta) = \tilde{g}(\theta) - \nabla_{\theta} L(\theta)$$

1. Magnitude <= YES! (e.g., Jastrzkebski et al. 2017)
2. Covariance structure <= YES! (e.g., Zhu et al. 2018)
3. Distribution class <= ? No!!! (this work)

Bernoulli? Gaussian? Levy?...

# Intuition

For quadratic loss, the generalization error

$$\mathbb{E}_{x, \theta_T} [\ell(x; \theta_T) - \ell(x; \theta_*)]$$

only depends on the first two moments of  $\theta_T$ , which only depend on the first two moments of  $v(\theta)$ .

$$\theta_{t+1} = \theta_t - \eta \underbrace{\nabla_{\theta} L(\theta_t)}_{\text{Linear}} - \eta v(\theta_t)$$

Noise matters! But noise class does not!!!

## A closer look at the noise of SGD

$$\underbrace{v_{\text{sgd}}(\theta)}_{\text{Gradient noise}} = \underbrace{\tilde{g}(\theta)}_{\nabla_{\theta} \mathcal{L}(\theta) \cdot \mathcal{W}_{\text{sgd}}} - \underbrace{\nabla_{\theta} L(\theta)}_{\nabla_{\theta} \mathcal{L}(\theta) \cdot \frac{1}{n} \mathbb{1}} = \nabla_{\theta} \mathcal{L}(\theta) \cdot \underbrace{\mathcal{V}_{\text{sgd}}}_{\text{Sampling noise}}$$

- Gradient matrix  $\nabla_{\theta} \mathcal{L}(\theta) = (\nabla_{\theta} \ell(x_1; \theta), \dots, \nabla_{\theta} \ell(x_n, \theta))$
- Sampling vector  $\mathcal{W}_{\text{sgd}} : \# \frac{1}{b} = b, \# 0 = n - b$
- Sampling noise  $\mathcal{V}_{\text{sgd}} = \mathcal{W}_{\text{sgd}} - \frac{1}{n}$

# Gradient noise vs. sampling noise

$$\underbrace{v(\theta)} = \underbrace{\nabla_{\theta} \mathcal{L}(\theta)} \cdot \underbrace{\mathcal{V}}$$

*Gradient noise*

*Gradient matrix*

*Sampling noise*

- State-dependent
- Noise of gradient

- State-dependent
- Deterministic

- **State-independent**
- Noise of mini-batch sampling

# Noisy gradient descent

$$\theta_{t+1} = \underbrace{\theta_t - \eta \nabla_{\theta} L(\theta_t)}_{\text{GD}} - \underbrace{\eta v(\theta)}_{\text{with noise}}$$

- in the same magnitude/covariance
- from different classes

Option 1: use gradient noise  $v(\theta)$  ☹️

Option 2: use sampling noise  $v(\theta) = \nabla_{\theta} \mathcal{L}(\theta) \cdot \mathcal{V}$  😊



# Multiplicative SGD (MSGD)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t) \cdot \mathcal{W}, \quad \mathcal{W} = \frac{1}{n} + \mathcal{V}$$

Algorithm:

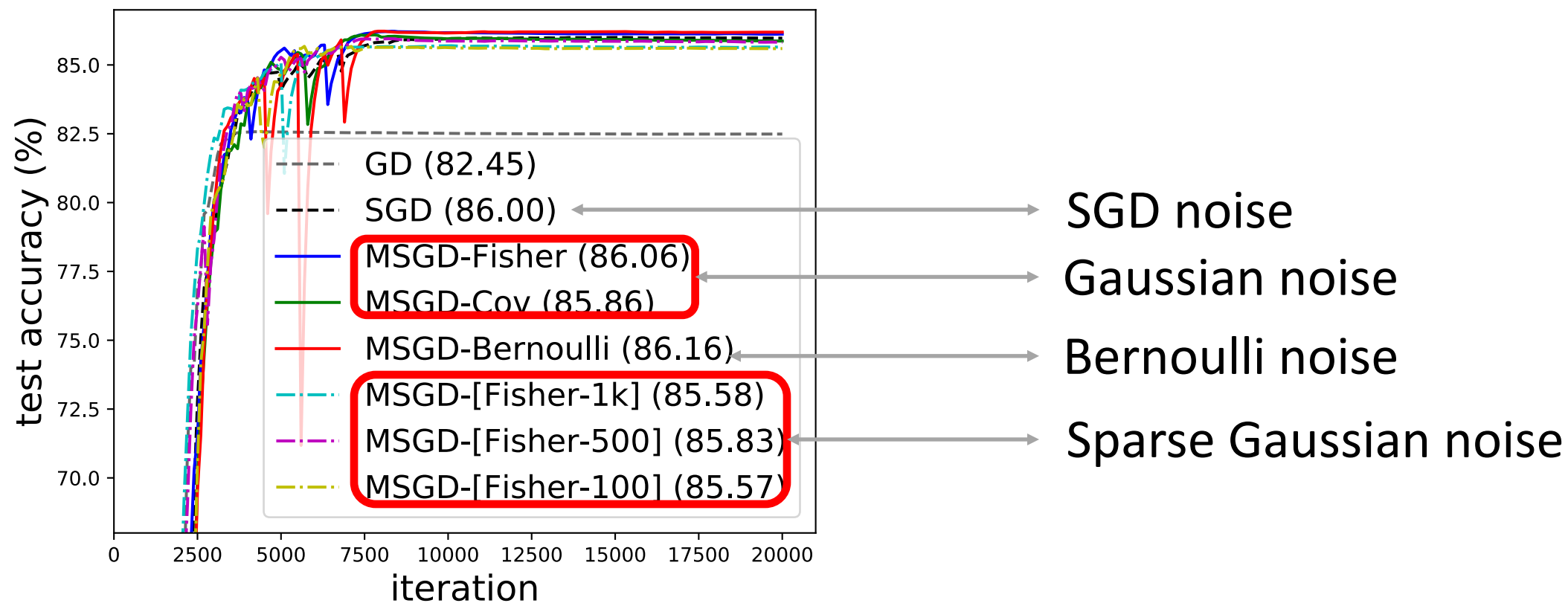
- |                                |   |
|--------------------------------|---|
| 1. Generate sampling vector    | $\mathcal{W} = 1/n + \mathcal{V}$                                   |
| 2. Compute randomized loss     | $\tilde{L}(\theta) = \mathcal{L}(\theta) \cdot \mathcal{W}$         |
| 3. Compute stochastic gradient | $\nabla_{\theta} \tilde{L}(\theta)$                                 |
| 4. Update parameters           | $\theta \leftarrow \theta - \eta \nabla_{\theta} \tilde{L}(\theta)$ |

# Injecting noise by MSGD

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t) \cdot \mathcal{W}, \quad \mathcal{W} = \frac{1}{n} + \mathcal{V}$$

- |                            |  |
|----------------------------|--|
| 1. SGD class               | $\mathcal{W}_{\text{sgd}} : \# \frac{1}{b} = b, \# 0 = n - b$  |
| 2. Gaussian class          | $\mathcal{W}_G \sim \mathcal{N}(1/n, \text{Var}[\mathcal{W}_{\text{sgd}}])$  |
| 3. “Bernoulli” class       | $\mathbb{P}\left(\mathcal{W}_B^{(i)} = \frac{1}{b}\right) = \frac{b}{n}, \mathbb{P}\left(\mathcal{W}_B^{(i)} = 0\right) = 1 - \frac{b}{n}$ |
| 4. “Sparse Gaussian” class | Mini-batch + Gaussian noise  |

# Experiments



Small SVHN. More results are available in the paper!

# Take Home



Get the paper!

1. Noise class is not crucial
2. Multiplicative SGD algorithm
3. Sampling noise perspective

Join our poster session for more details!