# Analyzing the Generalization Capability of SGLD Using Properties of Gaussian Channels

**Hao Wang**, Yizhe Huang, Rui Gao, Flavio P. Calmon

hao_wang@g.harvard.edu

yizhehuang@utexas.edu

rui.gao@mccombs.utexas.edu

flavio@seas.harvard.edu

**Harvard** John A. Paulson
**School of Engineering**
and Applied Sciences

TEXAS McCombs
The University of Texas at Austin
McCombs School of Business

# Outline

- Preliminary

  o Generalization analysis

  o An information-theoretic framework

- SGLD generalization bound

  o Definition

  o Our generalization bound

  o Experiments

- Generalization amplification by iteration

  o DP-SGD algorithm

  o Our generalization bound

- Related works and open questions

# Population risk minimization

Consider the following (non-convex) optimization problem:

$$\min_{w \in \mathcal{W}} \; L_\mu(w) \triangleq \mathbb{E}_{\mathsf{Z} \sim \mu} \left[ \ell(w, \mathsf{Z}) \right]$$

# Population risk minimization

Consider the following (non-convex) optimization problem:

$$\min_{w \in \mathcal{W}} \; L_\mu(w) \triangleq \mathbb{E}_{\mathsf{Z} \sim \mu}\left[\ell(w, \mathsf{Z})\right]$$

population risk

parameter

hypothesis class

# Population risk minimization

Consider the following (non-convex) optimization problem:

$$\min_{w \in \mathcal{W}} \; L_\mu(w) \triangleq \mathbb{E}_{\mathsf{Z} \sim \mu} \left[ \ell(w, \mathsf{Z}) \right]$$

loss function

data point
following distribution $\mu$

# Population risk minimization: Neural network

Consider the following (non-convex) optimization problem:

$$\min_{w \in \mathcal{W}} \; L_\mu(w) \triangleq \mathbb{E}_{\mathsf{Z} \sim \mu} \left[ \ell(w, \mathsf{Z}) \right]$$

weight matrices

$(f_w(\mathsf{X}) - \mathsf{Y})^2$ where $\mathsf{Z} = (\mathsf{X}, \mathsf{Y})$
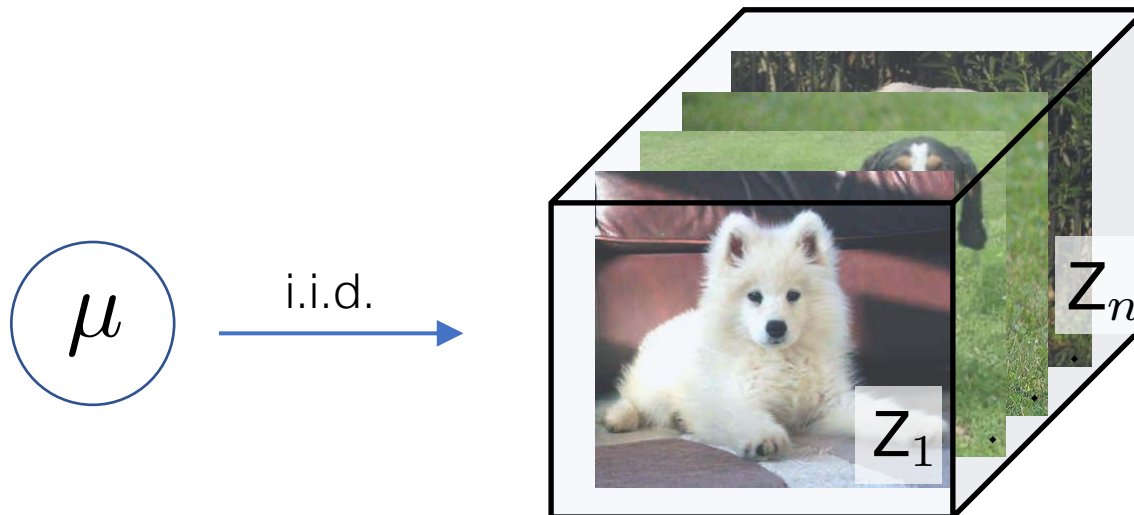
$f_w$

# Empirical risk minimization

Consider the following (non-convex) optimization problem:

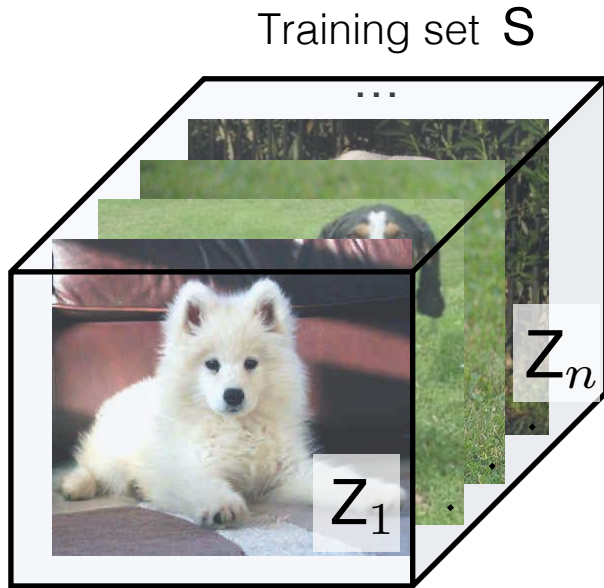$$\min_{w \in \mathcal{W}} \; L_\mu(w) \triangleq \mathbb{E}_{\mathsf{Z} \sim \mu} \left[ \ell(w, \mathsf{Z}) \right]$$

In practice,

$$\min_{w \in \mathcal{W}} \; L_{\mathsf{S}}(w) \triangleq \frac{1}{n} \sum_{i=1}^{n} \ell(w, \mathsf{Z}_i)$$

Training set S



$\mu$

i.i.d.

$\mathsf{Z}_1$    $\mathsf{Z}_n$

# Error decomposition

Training set S

$$\min_{w \in \mathcal{W}} \; L_{\mathsf{S}}(w)$$

output W

$Z_n$

$Z_1$

$$L_{\mu}(\mathsf{W}) = L_{\mathsf{S}}(\mathsf{W}) + \underline{L_{\mu}(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W})}$$

generalization gap

# Expected generalization gap

Training set S

$$\min_{w \in \mathcal{W}} L_S(w)$$

output W



$$\mathbb{E}\left[L_\mu(\mathsf{W}) - L_S(\mathsf{W})\right]$$

Expected generalization gap

# Today's focus

- VC-dimension

- Rademacher complexity

- Algorithmic stability

- PAC-Bayes

- Information theory

- …

# Information-theoretic generalization bound

**Theorem** (Xu and Raginsky, 2017). *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*
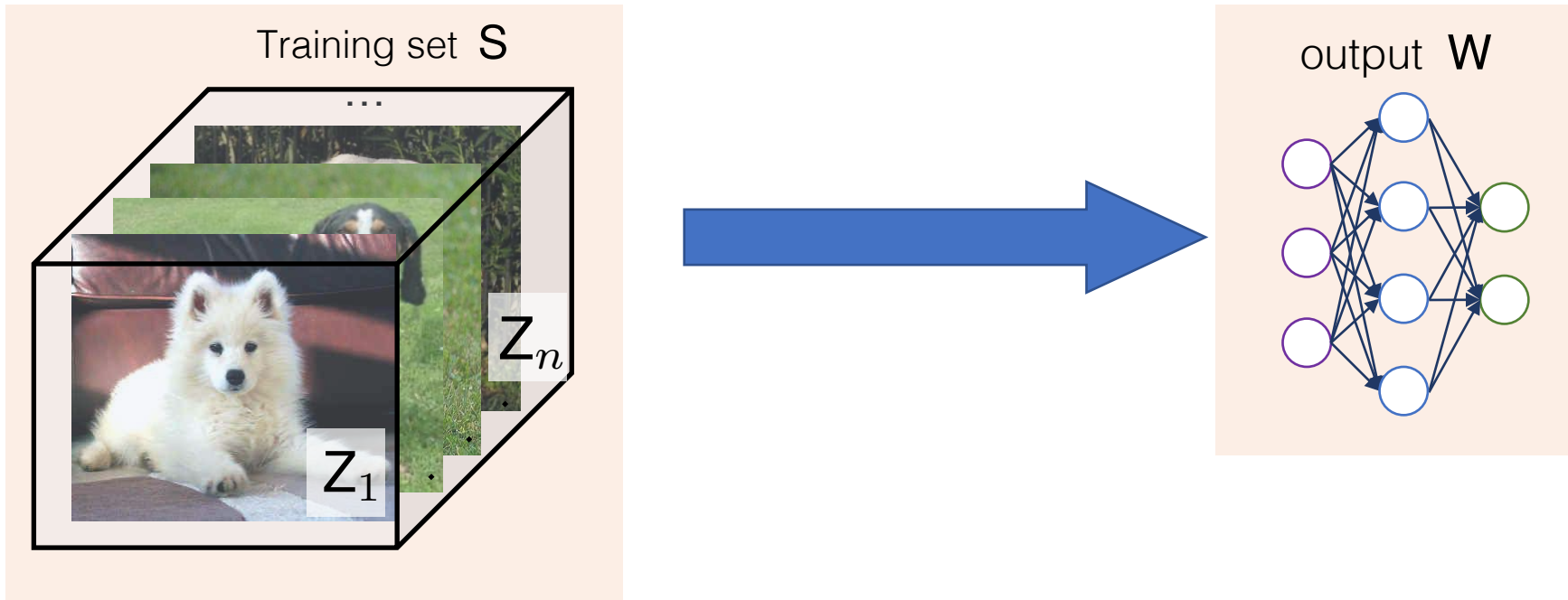
$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \sqrt{\frac{2\sigma^2}{n}I(\mathsf{W};\mathsf{S})}.$$

# Information-theoretic generalization bound

**Proposition** (Bu et al., 2020). *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{1}{n} \sum_{i=1}^{n} \sqrt{2\sigma^2 I(\mathsf{W}; \mathsf{Z}_i)}.$$



Training set $\mathsf{S}$

$\mathsf{Z}_i$

output $\mathsf{W}$

# Pros and cons

**Proposition** (Bu et al., 2020). *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{1}{n}\sum_{i=1}^{n}\sqrt{2\sigma^2 I(\mathsf{W}; \mathsf{Z}_i)}.$$

- Algorithm/distribution dependent
- Mild assumption
- ...

# Pros and cons

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{1}{n} \sum_{i=1}^{n} \sqrt{2\sigma^2 I(\mathsf{W}; \mathsf{Z}_i)}.$$

- Algorithm/distribution dependent

- Mild assumption

- …

A bounded loss is sufficient:
if $\ell(\cdot, \cdot) \in [0, A]$, then $\sigma = \frac{A}{2}$

# Pros and cons

**Proposition** (Bu et al., 2020). *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W}) \right] \right| \leq \frac{1}{n} \sum_{i=1}^{n} \sqrt{2\sigma^2 I(\mathsf{W}; \mathsf{Z}_i)}.$$

- Mutual information is hard to compute

# Outline

- Preliminary

  o Generalization analysis

  o An information-theoretic framework

- SGLD generalization bound

  o Definition

  o Our generalization bound

  o Experiments

- Generalization amplification by iteration

  o DP-SGD algorithm

  o Our generalization bound

- Related works and open questions

# SGLD is popular in practice

- Privacy guarantee

- Mitigate overfitting

- Easy to analyze in theory

- …

## tensorflow/**privacy**

Library for training machine learning models with privacy for training data

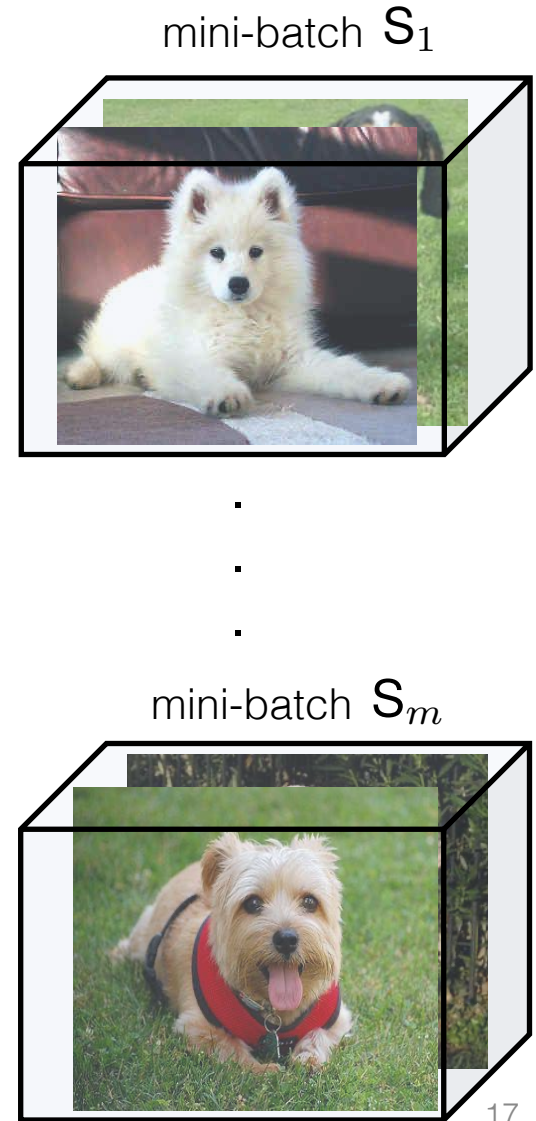| 🧑 42 | ⊙ 56 | ☆ 1k | ⅛ 321 | |
|---|---|---|---|---|
| Contributors | Issues | Stars | Forks | |

## Opacus

# SGLD: mini-batches

Training set $\mathsf{S}$



$\mathsf{z}_n$

$\mathsf{z}_1$

mini-batch $\mathsf{S}_1$

mini-batch $\mathsf{S}_m$

# SGLD: update rule

choose $W_0$ arbitrarily
for $t = 1, \cdots, T$

$$W_t = W_{t-1} - \eta_t \nabla_w \ell(W_{t-1}, S_{B_t}) + \sqrt{\frac{2\eta_t}{\beta_t}} N$$

learning rate

mini-batch gradient

mini-batch $S_1$



mini-batch $S_m$



18

[Welling and Teh, 2011]

# SGLD: update rule

choose $\mathsf{W}_0$ arbitrarily
for $t = 1, \cdots, T$

$$\mathsf{W}_t = \mathsf{W}_{t-1} - \eta_t \nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_{B_t}) + \sqrt{\frac{2\eta_t}{\beta_t}} \mathsf{N}$$

Gaussian
noise

inverse
temperature

# SGLD: output
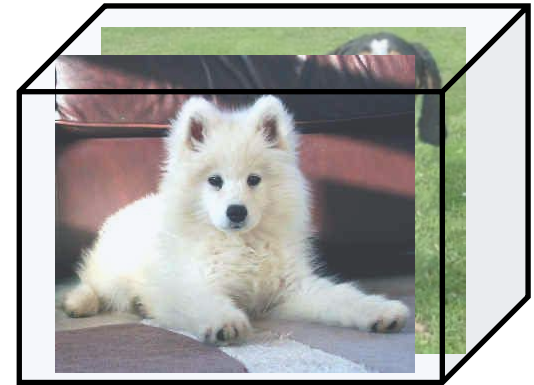
choose $W_0$ arbitrarily

for $t = 1, \cdots, T$

$$W_t = W_{t-1} - \eta_t \nabla_w \ell(W_{t-1}, S_{B_t}) + \sqrt{\frac{2\eta_t}{\beta_t}} N$$

output: $W = f(W_1, \cdots, W_T)$

# SGLD: output

choose $W_0$ arbitrarily
for $t = 1, \cdots, T$

$$W_t = W_{t-1} - \eta_t \nabla_w \ell(W_{t-1}, S_{B_t}) + \sqrt{\frac{2\eta_t}{\beta_t}} N$$

output: $W = f(W_1, \cdots, W_T)$

Examples:

- $f(W_1, \cdots, W_T) = W_T$

- $f(W_1, \cdots, W_T) = \frac{W_1 + \cdots + W_T}{T}$

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W}) \right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left( \nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j) \right)}.$$

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W}) \right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \text{Var}\left( \nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j) \right)}.$$

- $b$: mini-batch size

- $\sigma$: sub-Gaussian constant

- $n$: number of samples

- $\beta_t$: inverse temperature

- $\eta_t$: learning rate

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W}) \right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

- $b$: mini-batch size

- $\sigma$: sub-Gaussian constant

- $n$: number of samples

- $\beta_t$: inverse temperature

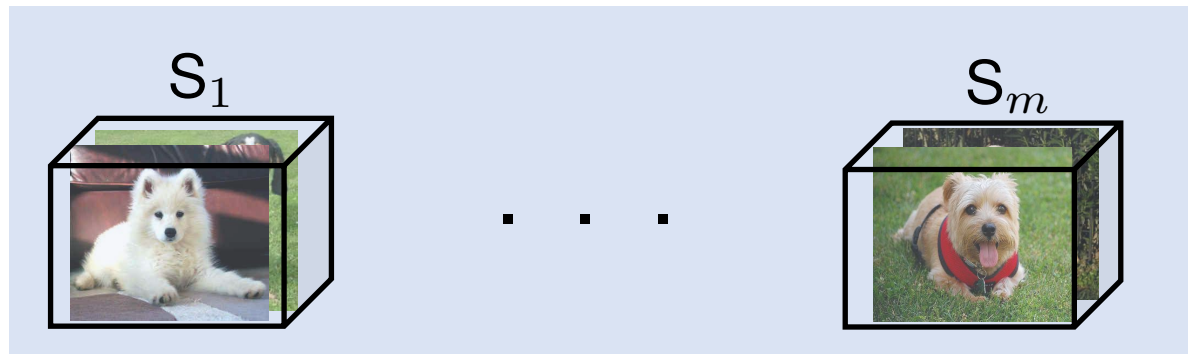- $\eta_t$: learning rate

Recall the update rule:
$$\mathsf{W}_t = \mathsf{W}_{t-1} - \eta_t \nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_{B_t}) + \sqrt{\frac{2\eta_t}{\beta_t}}\mathsf{N}$$

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$



mini-batch  $\mathsf{S}_1$  . . .  $\mathsf{S}_m$

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \textit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$
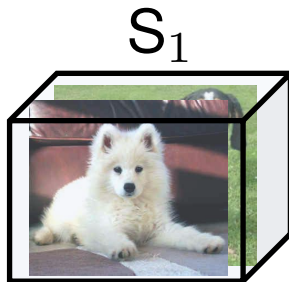
$\mathcal{T}_j$ contains the indices of iterations in which $\mathsf{S}_j$ is used

$\mathsf{S}_1$

mini-batch

$\cdot$ $\cdot$ $\cdot$

$\mathsf{S}_m$

# If the number of iterations increases,

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

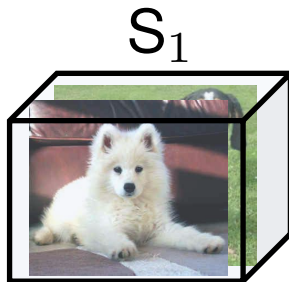$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W}) \right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \textit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

more terms

invariant

mini-batch

$\mathsf{S}_1$

$\cdots$

$\mathsf{S}_m$

27

# Our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]\right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

Measure sharpness of loss landscape

# Our generalization bound

**Theorem.** *Suppose $\ell(w, Z)$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*
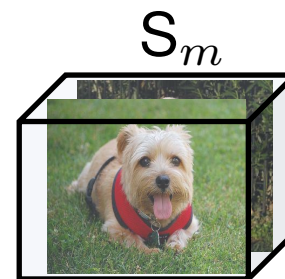
$$|\mathbb{E}[L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W})]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j))}.$$
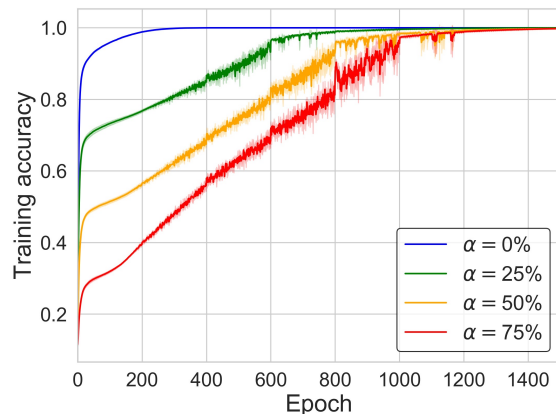
Variance of gradients

The journey matters more than the destination.

# Experiment on MNIST: Label corruption

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W}) \right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$
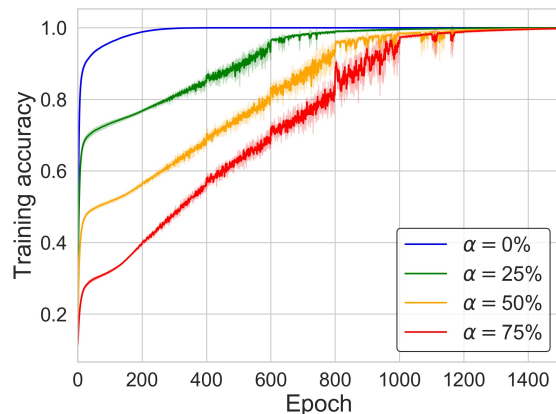
- Vary label corruption level (alpha)

- Train 3-layer neural networks on MNIST using SGLD

# Experiment on MNIST: Label corruption

**Theorem.** *Suppose $\ell(w, Z)$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W})\right] \right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$
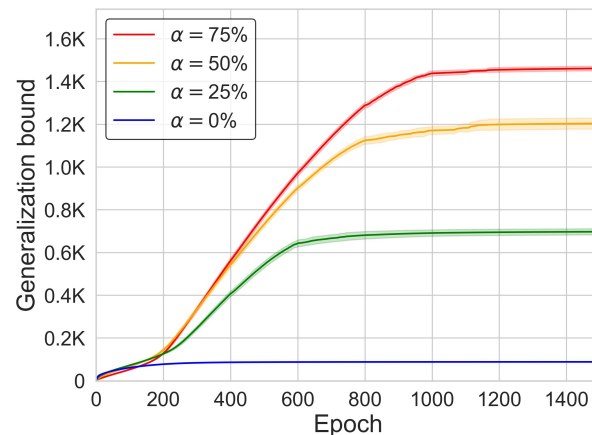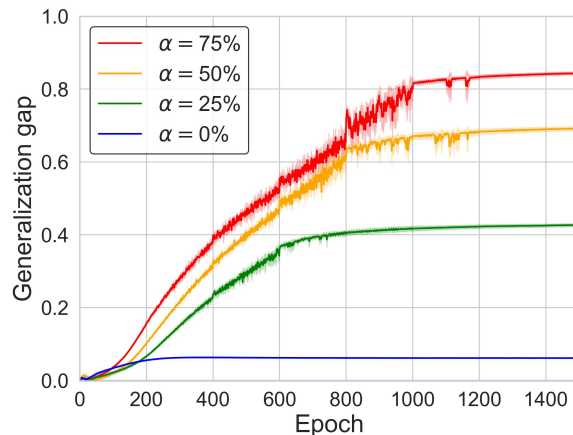


- Vary label corruption level (alpha)

- Train 3-layer neural networks on MNIST using SGLD

- Run 1500 epochs until the training accuracy = 1

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

# Experiment on MNIST: Label corruption

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]\right| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$
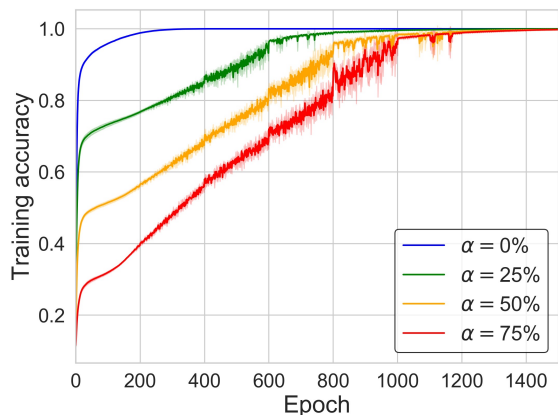
Both become stable

loss landscape



flat

When the algorithm converges, the variance of gradients become negligible.

# Experiment on MNIST: Network width

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\left| \mathbb{E}\left[ L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W}) \right] \right| \le \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot Var\left( \nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j) \right)}.$$
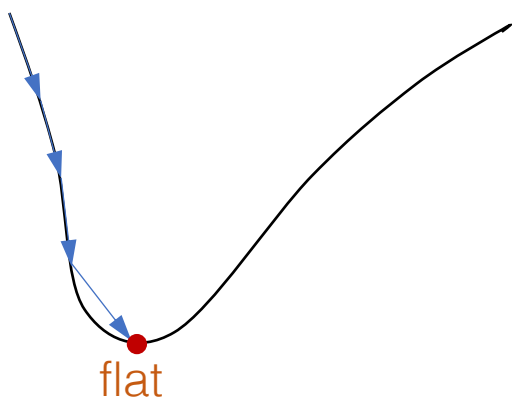
# Take away



Our bound can:
- be computed from data
- explain some generalization phenomena

# Outline

- Preliminary

  o    Generalization analysis

  o    An information-theoretic framework

- SGLD generalization bound

  o    Definition

  o    Our generalization bound

  o    Experiments

- Generalization amplification by iteration

  o    DP-SGD algorithm

  o    Our generalization bound

- Related works and open questions

# Recall our generalization bound

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_{\mathsf{S}}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \textit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

Here $\mathsf{W} = f(\mathsf{W}_1, \cdots, \mathsf{W}_T)$ for an arbitrary function $f$

# Upside

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$
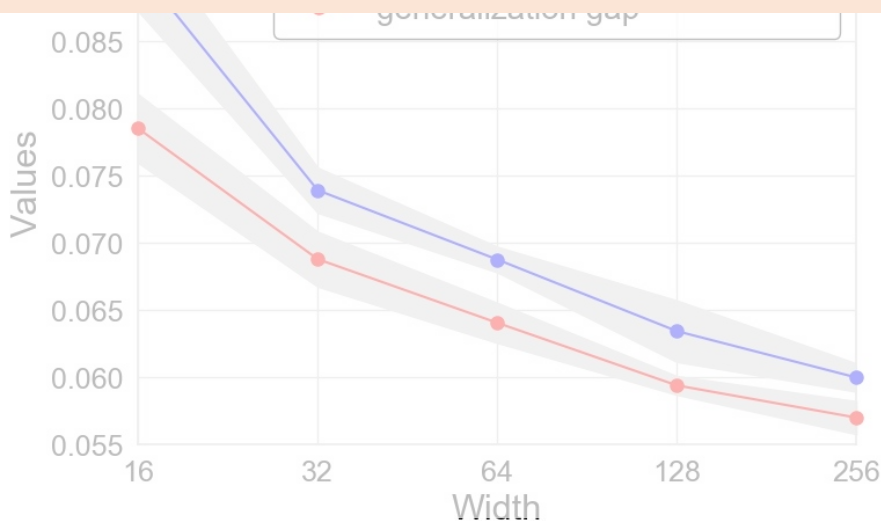
Here $\mathsf{W} = f(\mathsf{W}_1, \cdots, \mathsf{W}_T)$ for an arbitrary function $f$

Can be applied to many settings

Examples:

- $f(\mathsf{W}_1, \cdots, \mathsf{W}_T) = \mathsf{W}_T$
- $f(\mathsf{W}_1, \cdots, \mathsf{W}_T) = \frac{\mathsf{W}_1 + \cdots + \mathsf{W}_T}{T}$

# Downside

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$\sup_{f} |\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

Here $\mathsf{W} = f(\mathsf{W}_1, \cdots, \mathsf{W}_T)$ for an arbitrary function $f$

A uniform bound for any function $f$

# Question

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}) - L_\mathsf{S}(\mathsf{W})\right]| \leq \frac{\sqrt{2b}\sigma}{2n} \sum_{j=1}^{m} \sqrt{\sum_{t \in \mathcal{T}_j} \beta_t \eta_t \cdot \mathit{Var}\left(\nabla_w \ell(\mathsf{W}_{t-1}, \mathsf{S}_j)\right)}.$$

Here $\mathsf{W} = f(\mathsf{W}_1, \cdots, \mathsf{W}_T)$ for an arbitrary function $f$

For $\mathsf{W} = \mathsf{W}_T$, can we have a sharper bound?

# Projected Differentially-Private SGD (DP-SGD)

choose $W_0$ arbitrarily

for $t = 1, \cdots, T$

$$W_t = \text{Proj}_{\mathcal{W}}\left(W_{t-1} - \eta\left(g(W_{t-1}, Z_t) + N\right)\right)$$

output: $W_T$

Assumptions:

- sampling without replacement

- $\|g(w, z)\|_2 \leq K$ for any $w, z$

# Our generalization bound: time-decaying factor

**Theorem.** *Suppose $\ell(w, \mathsf{Z})$ is $\sigma$-sub-Gaussian under $\mu$ for all $w$.*

$$|\mathbb{E}\left[L_\mu(\mathsf{W}_T) - L_{\mathsf{S}}(\mathsf{W}_T)\right]| \leq \frac{2\sigma}{n} \sum_{t=1}^{T} \sqrt{\mathit{Var}\left(g\left(\mathsf{W}_{t-1}, \mathsf{Z}\right)\right) \cdot q^{T-t}}.$$

$$q \in (0, 1)$$

Enables the impact of early iterations to reduce with time

# Proof

Key proof techniques: properties of Gaussian channels

Step 1:  $|\mathbb{E}\left[L_\mu(\mathbf{W}_T) - L_S(\mathbf{W}_T)\right]| \leq \frac{\sqrt{2}\sigma}{n} \sum \sqrt{I(\mathbf{W}_T; \mathbf{Z}_t)}$       [Bu et al., 2020]
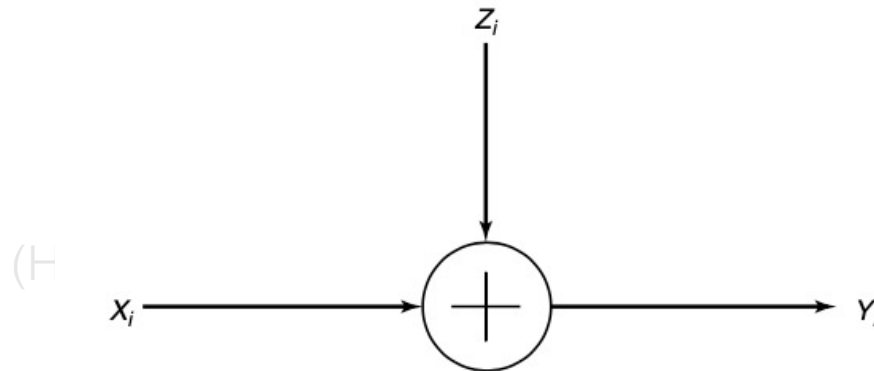
Step 2:

(⊢



**FIGURE 9.1.** Gaussian channel.

Step 3:

(How we obtain a computable bound)

43

Image is from [Cover and Thomas, 1999]

# Outline

- **Preliminary**
  - Generalization analysis
  - An information-theoretic framework

- **SGLD generalization bound**
  - Definition
  - Our generalization bound
  - Experiments

- **Generalization amplification by iteration**
  - DP-SGD algorithm
  - Our generalization bound

- **Related works and open questions**

# Related works and open questions

- SGLD generalization bounds: [Mou et al., 2018], [Li et al., 2019], [Pensia et al., 2018], [Bu et al., 2020], [Negrea et al., 2019], [Haghifam et al., 2020], [Rodriguez-Galvez et al., 2020] [Neu, 2021]

- Privacy amplification by iteration: [Feldman et al., 2018], [Asoodeh et al., 2020]

# Related works and open questions

- SGLD generalization bounds: [Mou et al., 2018], [Li et al., 2019], [Pensia et al., 2018], [Bu et al., 2020], [Negrea et al., 2019], [Haghifam et al., 2020], [Rodriguez-Galvez et al., 2020] [Neu, 2021]

- Privacy amplification by iteration: [Feldman et al., 2018], [Asoodeh et al., 2020]

- Privacy can also be amplified by subsampling and shuffling. Can they improve algorithmic generalization?

- New definition of sharpness

- Tighten the decay factor.

[Chaudhuri and Mishra, 2006; Erlingsson et al., 2019]

# Related works and open questions

- SGLD generalization bounds: [Mou et al., 2018], [Li et al., 2019], [Pensia et al., 2018], [Bu et al., 2020], [Negrea et al., 2019], [Haghifam et al., 2020], [Rodriguez-Galvez et al., 2020] [Neu, 2021]

- Privacy amplification by iteration: [Feldman et al., 2018], [Asoodeh et al., 2020]

- Privacy can also be amplified by subsampling and shuffling. Can they improve algorithmic generalization?

- New definition of sharpness

- Tighten the decay factor.

# Related works and open questions

- SGLD generalization bounds: [Mou et al., 2018], [Li et al., 2019], [Pensia et al., 2018], [Bu et al., 2020], [Negrea et al., 2019], [Haghifam et al., 2020], [Rodriguez-Galvez et al., 2020] [Neu, 2021]

- Privacy amplification by iteration: [Feldman et al., 2018], [Asoodeh et al., 2020]

- Privacy can also be amplified by subsampling and shuffling. Can they improve algorithmic generalization?

- New definition of sharpness

- Tighten the decay factor.

# Thanks for watching!