

# Generalization of Neural Network

Qingru Zhang

Apex Data & Knowledge Management Lab  
Shanghai Jiao Tong University



# Content

- **Generalization Bound**
- **PAC-Bayesian Framework**
- **Randomization Test**
- **Overparameterization**
- **Conclusion**

# Generalization Bound

## Definition:

$\ell(y, y')$  - loss function

$\mathcal{H}$  - hypothesis space

$h(x)$  - prediction of hypothesis  $h \in \mathcal{H}$  on sample  $x$

$L(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, h(x))]$  - expected loss of  $h$

$\hat{L}(h) = \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(x_i))$  - empirical loss of  $h$

generalization error =  $L(h) - \hat{L}(h)$

# Generalization Bound

## Generalization Bound Overview:

With probability  $1 - \delta$

$$\hat{L}(h) \leq L(h) \leq \hat{L}(h) + \Omega(m, R, \delta)$$

$E_{train}$        $E_{test}$

Smaller  $\delta$ , larger  $\Omega$

$m$  is the number of training data  $\longrightarrow$  Larger  $m$ , smaller  $\Omega$

$R$  is the “capacity” of your model  $\longrightarrow$  Larger  $R$ , larger  $\Omega$   
(“size” of the function set)

## Generalization Bound Framework:

- VC-dimension
- PAC-Bayesian
- Rademacher complexity

$$\hat{\mathfrak{R}}_n(\mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i) \right]$$

# Generalization Bound

## Current Framework:

- VC-dimension

$$\text{VC-dim} = \tilde{O}(d * \dim(\mathbf{w}))$$

- PAC-Bayesian Framework
- Rademacher complexity

$$\hat{\mathfrak{R}}_n(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i) \right]$$

# General PAC-Bayesian Theorem

## Preliminaries:

- $P$  is a prior distribution on  $h$
- $Q$  is a posterior distribution on  $h$
- Take the loss expectation on  $Q$ .

### Gibbs Risk / Linear Loss

The stochastic Gibbs classifier  $G_Q(x)$  draws  $h' \in \mathcal{H}$  according to  $Q$  and output  $h'(x)$ .

$$\begin{aligned} R_D(G_Q) &= \mathbf{E}_{(x,y) \sim D} \mathbf{E}_{h \sim Q} \mathbf{I}[h(x) \neq y] \\ &= \mathbf{E}_{h \sim Q} \mathcal{L}_D^{\ell_{01}}(h), \end{aligned}$$

where  $\ell_{01}(h, x, y) = \mathbf{I}[h(x) \neq y]$ .

# General PAC-Bayesian Theorem

$\Delta$ -function: “distance” between  $\hat{R}_S(G_Q)$  et  $R_D(G_Q)$

Convex function  $\Delta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ .

General theorem (Bégin et al. (2014b, 2016); Germain (2015))

*For any distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ , for any set  $\mathcal{H}$  of voters, for any distribution  $P$  on  $\mathcal{H}$ , for any  $\delta \in (0, 1]$ , and for any  $\Delta$ -function, we have, with probability at least  $1 - \delta$  over the choice of  $S \sim D^m$ ,*

$$\forall Q \text{ on } \mathcal{H} : \quad \Delta\left(\hat{R}_S(G_Q), R_D(G_Q)\right) \leq \frac{1}{m} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(m)}{\delta} \right],$$

where

$$\mathcal{I}_\Delta(m) = \sup_{r \in [0, 1]} \left[ \sum_{k=0}^m \underbrace{\binom{m}{k} r^k (1-r)^{m-k}}_{\text{Bin}(k; m, r)} e^{m \Delta(\frac{k}{m}, r)} \right].$$

# General PAC-Bayesian Theorem

$$\Pr_{S \sim D^m} \left( \forall Q \text{ on } \mathcal{H} : \Delta \left( \hat{R}_S(G_Q), R_D(G_Q) \right) \leq \frac{1}{m} \left[ \text{KL}(Q \| P) + \ln \frac{\mathcal{I}_\Delta(m)}{\delta} \right] \right) \geq 1 - \delta. \quad \text{Proof.}$$

$$m \cdot \Delta \left( \mathbf{E}_{h \sim Q} \hat{\mathcal{L}}_S^\ell(h), \mathbf{E}_{h \sim Q} \mathcal{L}_D^\ell(h) \right)$$

Jensen's Inequality

$$\leq \mathbf{E}_{h \sim Q} m \cdot \Delta \left( \hat{\mathcal{L}}_S^\ell(h), \mathcal{L}_D^\ell(h) \right)$$

Change of measure

$$\leq \text{KL}(Q \| P) + \ln \mathbf{E}_{h \sim P} e^{m \Delta \left( \hat{\mathcal{L}}_S^\ell(h), \mathcal{L}_D^\ell(h) \right)}$$

Markov's Inequality

$$\leq_{1-\delta} \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{S' \sim D^m} \mathbf{E}_{h \sim P} e^{m \cdot \Delta \left( \hat{\mathcal{L}}_{S'}^\ell(h), \mathcal{L}_D^\ell(h) \right)}$$

Expectation swap

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \mathbf{E}_{S' \sim D^m} e^{m \cdot \Delta \left( \hat{\mathcal{L}}_{S'}^\ell(h), \mathcal{L}_D^\ell(h) \right)}$$

Binomial law

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathbf{E}_{h \sim P} \sum_{k=0}^m \text{Bin}(k; m, \mathcal{L}_D^\ell(h)) e^{m \cdot \Delta \left( \frac{k}{m}, \mathcal{L}_D^\ell(h) \right)}$$

Supremum over risk

$$\leq \text{KL}(Q \| P) + \ln \frac{1}{\delta} \sup_{r \in [0,1]} \left[ \sum_{k=0}^m \text{Bin}(k; m, r) e^{m \Delta \left( \frac{k}{m}, r \right)} \right]$$

$$= \text{KL}(Q \| P) + \ln \frac{1}{\delta} \mathcal{I}_\Delta(m).$$

□



# Norm-based PAC-Bayesian Bound

## Margin loss

$$L_\gamma(f_{\mathbf{w}}) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ f_{\mathbf{w}}(\mathbf{x})[y] \leq \gamma + \max_{j \neq y} f_{\mathbf{w}}(\mathbf{x})[j] \right]$$

## PAC-Bayesian Lemma

$$L_0(f_{\mathbf{w}}) \leq \hat{L}_\gamma(f_{\mathbf{w}}) + 4 \sqrt{\frac{KL(\mathbf{w} + \mathbf{u} \| P) + \ln \frac{6m}{\delta}}{m-1}}.$$

## Perturbation Bound

$$|f_{\mathbf{w}+\mathbf{u}}(\mathbf{x}) - f_{\mathbf{w}}(\mathbf{x})|_2 \leq eB \left( \prod_{i=1}^d \|W_i\|_2 \right) \sum_{i=1}^d \frac{\|U_i\|_2}{\|W_i\|_2}.$$

## Generalization Bound

$$L_0(f_{\mathbf{w}}) \leq \hat{L}_\gamma(f_{\mathbf{w}}) + \mathcal{O} \left( \sqrt{\frac{B^2 d^2 h \ln(dh) \prod_{i=1}^d \|W_i\|_2^2 \sum_{i=1}^d \frac{\|W_i\|_F^2}{\|W_i\|_2^2} + \ln \frac{dm}{\delta}}{\gamma^2 m}} \right).$$

# No-vacuous PAC-Bayesian Bound

**A compression approach:**

- **different network has different compressibility**
- **Construct prior P according to compressibility**

$$\pi_c(h) = \frac{1}{Z} m(|h|_c) 2^{-|h|_c}, \text{ where } Z = \sum_{h \in \mathcal{H}_c} m(|h|_c) 2^{-|h|_c}.$$

- **Give larger probability to more compressible network**
- **Reduce KL-distance between P and Q**

$$\begin{aligned} \text{KL}(\rho_{S,C,Q}, \pi) &\leq (k \lceil \log r \rceil + |S|_c + |C|_c) \log 2 - \log m(k \lceil \log r \rceil + |S|_c + |C|_c) \\ &\quad + \sum_{i=1}^k \text{KL} \left( \text{Normal}(c_{q_i}, \sigma^2), \sum_{j=1}^r \text{Normal}(c_j, \tau^2) \right). \end{aligned}$$

# No-vacuous PAC-Bayesian Bound

## Bound:

$$\text{KL}(\rho_{S,C,Q}, \pi) \leq (k \lceil \log r \rceil + |S|_c + |C|_c) \log 2 - \log m(k \lceil \log r \rceil + |S|_c + |C|_c) \\ + \sum_{i=1}^k \text{KL}\left(\text{Normal}(c_{q_i}, \sigma^2), \sum_{j=1}^r \text{Normal}(c_j, \tau^2)\right).$$

## Experiments:

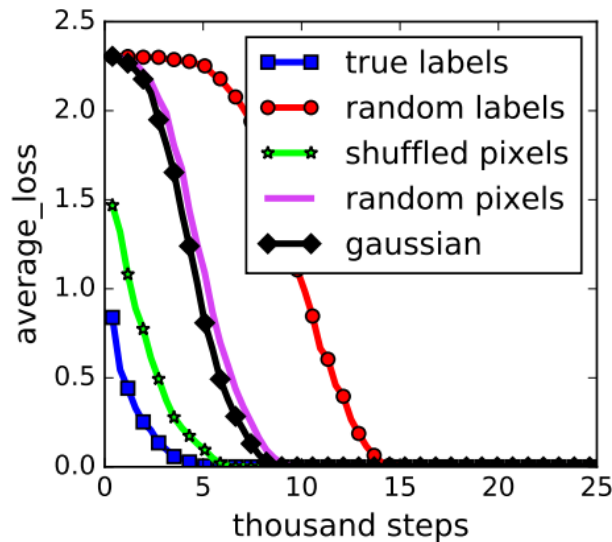
Dataset	Orig. size	Comp. size	Robust. Adj.	Eff. Size	Error Bound	
					Top 1	Top 5
MNIST	168.4 KiB	8.1 KiB	1.88 KiB	6.23 KiB	< 46 %	NA
ImageNet	5.93 MiB	452 KiB	102 KiB	350 KiB	< 96.5 %	< 89 %

# Randomization tests

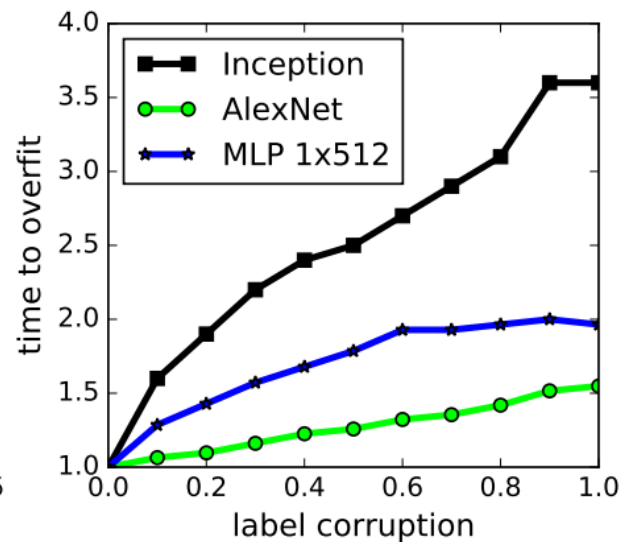
**Test Insight:** randomizing labels alone to force the generalization error jumping up without changing the model.

**Result:**

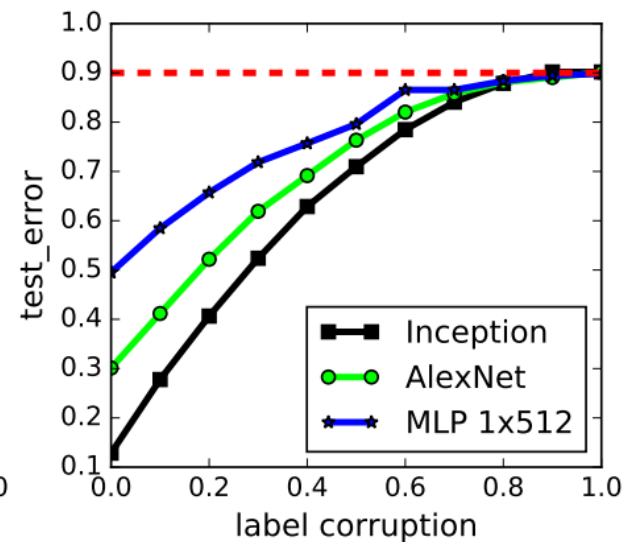
*Deep neural networks easily fit random labels.*



(a) learning curves



(b) convergence slowdown



(c) generalization error growth

# Randomization tests

## Expectation:

- Since initially the label is uncorrelated
- Large predictions errors are back-propagated
- Make large gradients for parameter updates

## Observation:

- No need to change the learning rate schedule
- Once the fitting starts, it converges quickly
- it converges to (over) fit the training set perfectly.

# Challenge for Complexity Measures

## Rademacher complexity

$$\hat{\mathfrak{R}}_n(\mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i) \right]$$

- Since many neural networks fit the training set with random labels perfectly, we expect that  $R \approx 1$ , which is a trivial upper bound.

## VC-dimension

$$\text{VC-dim} = \tilde{O}(d * \dim(\mathbf{w}))$$

- When the number of parameters is more than the number of samples, it becomes too weak.

# The role of regularization

## Regularization technique:

- Data augmentation, Weight Decay, Dropout.

## Result:

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
		(fitting random labels)	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
		(fitting random labels)	no	100.0	10.12
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
		(fitting random labels)	no	99.82	9.86
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
		(fitting random labels)	no	100.0	10.48
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
		(fitting random labels)	no	99.34	10.61

# The role of regularization

## Conclusion:

**Both regularization techniques help to improve the generalization performance**

**But turned off all regularizers, all of the models still generalize very well.**

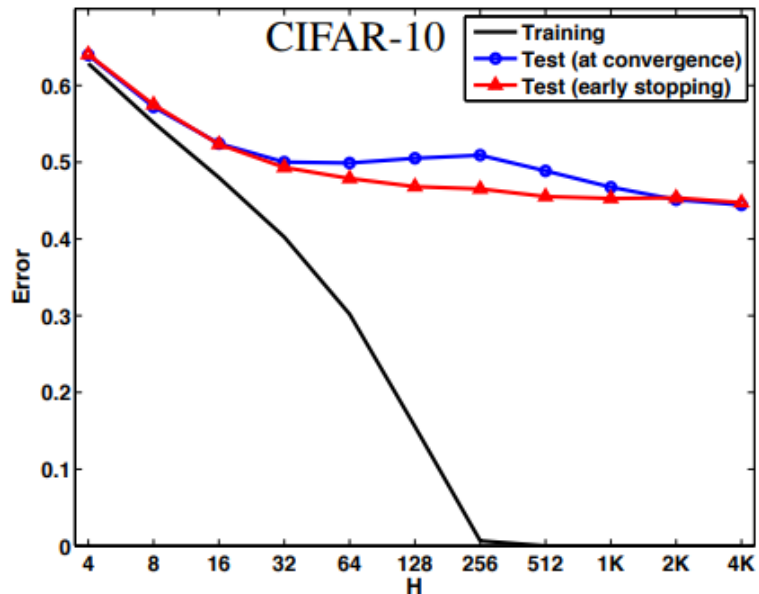
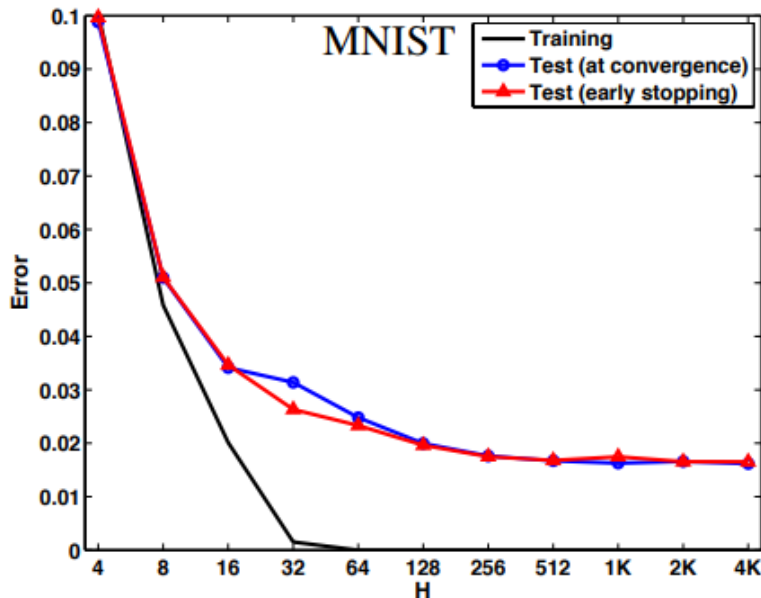
## Some Extension:

- **SGD acts as an implicit regularizer**
- **For linear models, SGD always converges to a solution with small norm.**
- **What the properties are inherited by models that were trained using SGD?**



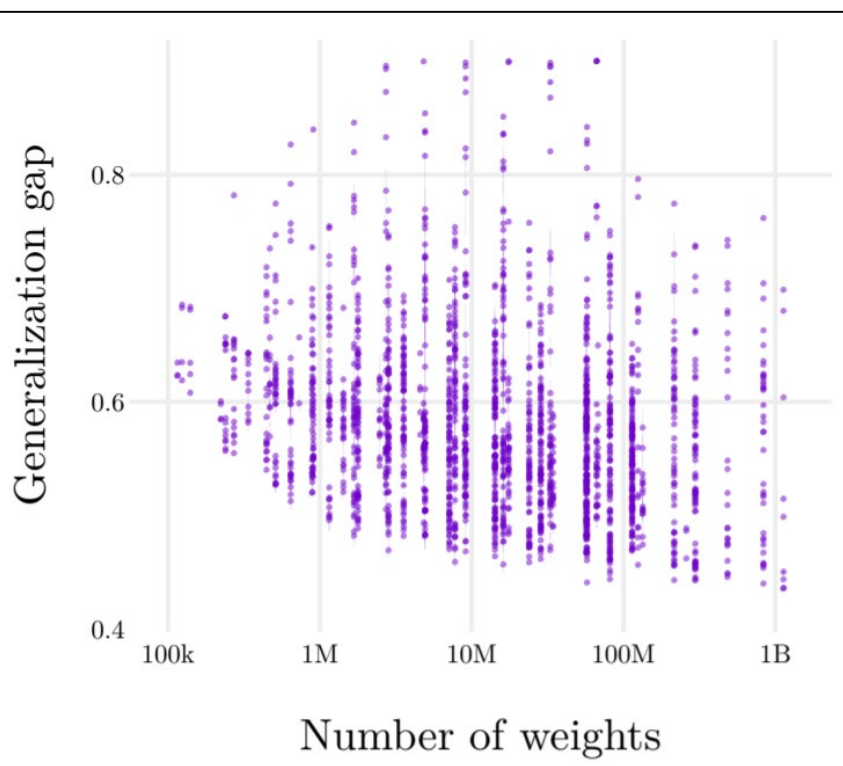
# Overparameterized Network

- Overparameterized Network still can generalize
- *IN SEARCH OF THE REAL INDUCTIVE BIAS*  
point out implicit regularization of SGD

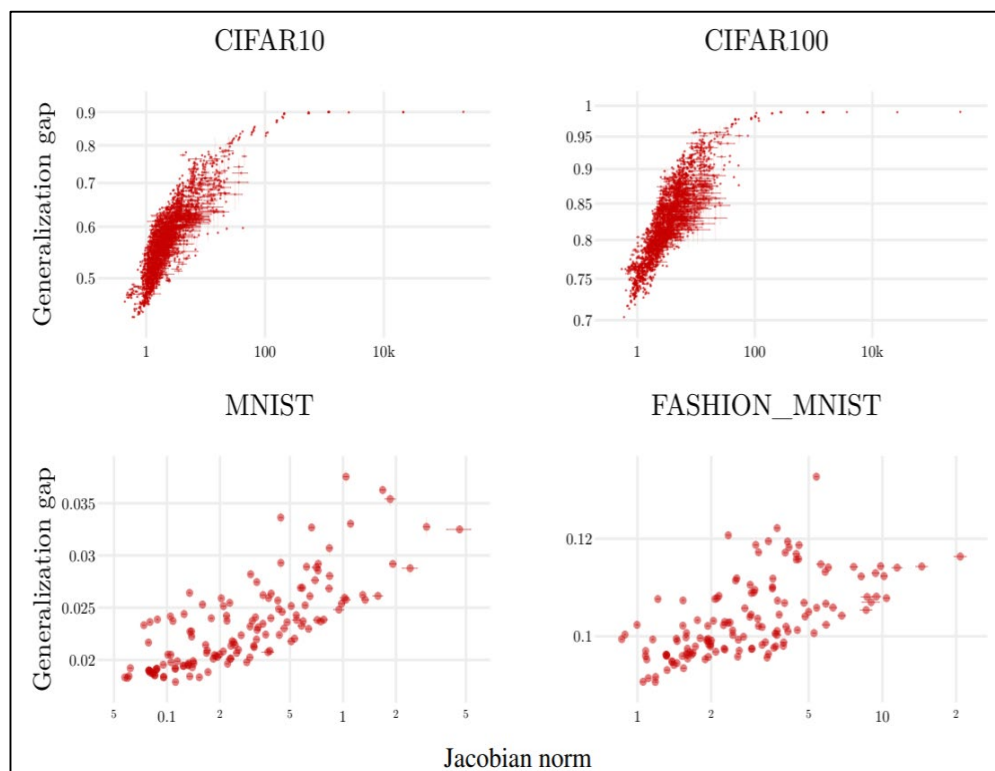


# Overparameterized Network

- Overparameterized Network still can generalize
- *SENSITIVITY AND GENERALIZATION IN NEURAL NETWORKS* focus on norm of Jacobian



CIFAR-10 100% training accuracy



F-norm

# Overparameterized Network

## Conclusion

- Overparameterized Network's capacity is large
- But it does not overfit and still can generalize
- The reason is unclear yet.
  - SGD's implicit regularization
  - Network regularizes itself
  - .....

# Conclusion

- **Overparameterized network still generalize**
- **Neural network can fit any dataset, random label or true label**
- **Theoretic analysis should consider dataset structure.**
- **Neural network owns some properties:**
  - **robust to perturbation of parameters**
  - **compressibility**
- **State-of-art theoretic bounds are still far from being effective**

# Q & A

# Thanks for Listening!