

Generalization of neural networks: A Hessian view (Some partial progress)

Hongyang Ryan Zhang

Assistant Professor, Khoury College of Computer Sciences
177 Huntington Ave (22nd floor), Northeastern University, Boston

December 4, 2024

Motivating questions

Some intriguing questions related to the working of GPT-3: How does a large model pretrained on a large corpus of unlabeled data learn, and why is it adaptable to multiple downstream tasks?

Motivating questions

Some intriguing questions related to the working of GPT-3: How does a large model pretrained on a large corpus of unlabeled data learn, and why is it adaptable to multiple downstream tasks?

- Related to research on the generalization of over-parameterized models (2017): Modern deep networks (e.g., ResNet, BERT) have more parameters than dataset labels to memorize training data [ZBH+21]

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

Oriol Vinyals
Google DeepMind
vinyals@google.com

ABSTRACT

Despite their massive size, successful deep artificial neural networks can exhibit a remarkably small difference between training and test performance. Conventional wisdom attributes small generalization error either to properties of the model family, or to the regularization techniques used during training.

Generalization of language models

Recently, there has been growing interest in the generalization capability of LLMs (broadly defined)

ACL 2025 Theme Track: Generalization of NLP Models

Following the success of the ACL 2020-2024 Theme tracks, we are happy to announce that ACL 2025 will have a new theme with the goal of reflecting and stimulating discussion about the current state of development of the field of NLP.

Generalization is crucial for ensuring that models behave robustly, reliably, and fairly when making predictions on data different from their training data. Achieving good generalization is critically important for models used in real-world applications, as they should emulate human-like behavior. Humans are known for their ability to generalize well, and models should aspire to this standard.

The theme track invites empirical and theoretical research and position and survey papers reflecting on the Generalization of NLP Models. The possible topics of discussion include (but are not limited to) the following:

- How can we enhance the generalization of NLP models across various dimensions—compositional, structural, cross-task, cross-lingual, cross-domain, and robustness?
- What factors affect the generalization of NLP models?
- What are the most effective methods for evaluating the generalization capabilities of NLP models?
- While Large Language Models (LLMs) significantly enhance the generalization of NLP models, what are the key limitations of LLMs in this regard?

Generalization of language models

Understanding what information is presented in a natural language and how we may quantify content is certainly not a new research quest:::

- The notion of entropy for measuring information in a language dates back to very classical work [[Sha48](#)]
- Applying entropy to N -gram statistics gives a measure of information (more recent instantiations include next-word prediction in GPT-3, and next-sentence prediction in BERT)

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x,y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- Rademacher complexity [BFT17]
- PAC-Bayes [AGN+18; JLZ22]: data-dependent bounds
- Neural tangent kernels [ADH+19]: high width, fixed random matrix
- Implicit regularization [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- Benign overfitting [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- PAC-Bayes [AGN+18; JLZ22]: data-dependent bounds
- Neural tangent kernels [ADH+19]: high width, fixed random matrix
- Implicit regularization [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- Benign overfitting [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Generalization of learning with IID samples

Training: given n samples from D , $(x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$,

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

Test: expected risk over a random sample of D

$$L(f_W) = \mathbb{E}_{(x;y) \sim D} [\ell(f_W(x); y)]$$

Generalization analysis of ML algorithms: given a hypothesis space of neural nets, how do we quantify generalization gap $L(f_W) - \hat{L}(f_W)$?

- **Rademacher complexity** [BFT17]
- **PAC-Bayes** [AGN+18; JLZ22]: data-dependent bounds
- **Neural tangent kernels** [ADH+19]: high width, fixed random matrix
- **Implicit regularization** [LMZ18]: starting from a small, random initialization, SGD searches inside a low-rank space
- **Benign overfitting** [BLL+20]

See a recent textbook by Tong Zhang [Zha23] (UIUC) for further references

Uniform convergence bounds on Lipschitz-continuous neural networks

Neural network setup

- f_W : A multi-layer neural net with weight matrices $W = [W_1; W_2; \dots; W_l]$
- D : An unknown distribution over feature space times label space $X \times Y$

Norm bounds for neural nets

Hypothesis space with bounded L_2 and Frobenius norms:

$$H = \left\{ \sum_{i=1}^n \sum_{j=1}^{r_1} x_j W_1 x_j + s_1 \sum_{i=1}^{r_2} x_j W_2 x_j + s_2 \sum_{i=1}^{r_3} x_j W_3 x_j + \dots + s_{l-1} \sum_{i=1}^{r_l} x_j W_l x_j + s_l \right\} \quad (1)$$

- For any l -layer deep net f_W whose weight matrices W belong to H , with Lipschitz-continuous activation functions at every layer, the following holds w.h.p.

$$L(f_W) \leq \hat{L}(f_W) \cdot \sqrt{\frac{\sum_{i=1}^l \frac{Q_i}{s_i^2} + \sum_{i=1}^l P_i r_i^2}{n}}$$

- The bound grows exponentially with depth, and is tight in the worst-case

Uniform convergence bounds on Lipschitz-continuous neural networks

Neural network setup

- f_W : A multi-layer neural net with weight matrices $W = [W_1; W_2; \dots; W_l]$
- D : An unknown distribution over feature space times label space $X \times Y$

Norm bounds for neural nets

- Hypothesis space with bounded ℓ_2 and Frobenius norms:

$$H = \left\{ \begin{matrix} \sum_{i=1}^n k W_1 k_2 \cdot s_{i1}; k W_2 k_2 \cdot s_{i2}; \dots; k W_l k_2 \cdot s_{il} \\ k W_1 k_F \cdot r_{i1}; k W_2 k_F \cdot r_{i2}; \dots; k W_l k_F \cdot r_{il} \end{matrix} \right\} \quad (1)$$

- For any l -layer deep net f_W whose weight matrices W belong to H , with Lipschitz-continuous activation functions at every layer, the following holds w.h.p.

$$L(f_W) \leq \hat{L}(f_W) \cdot \sqrt{\frac{\sum_{i=1}^l (Q_i \sum_{j=1}^n s_j^2 + P_i \sum_{j=1}^n r_j^2)}{n}}$$

- The bound grows exponentially with depth, and is tight in the worst-case

Uniform convergence bounds on Lipschitz-continuous neural networks

Neural network setup

- f_W : A multi-layer neural net with weight matrices $W = [W_1; W_2; \dots; W_l]$
- D : An unknown distribution over feature space times label space $X \times Y$

Norm bounds for neural nets

- Hypothesis space with bounded ℓ_2 and Frobenius norms:

$$H = \left\{ \begin{array}{l} \sum_{i=1}^n k W_1 k_2 \quad s_1; k W_2 k_2 \quad s_2; \dots; k W_l k_2 \quad s_l; \\ k W_1 k_F \quad s_1 r_1; k W_2 k_F \quad s_2 r_2; \dots; k W_l k_F \quad s_l r_l \end{array} \right\} \quad (1)$$

- For any l -layer deep net f_W whose weight matrices W belong to H , with Lipschitz-continuous activation functions at every layer, the following holds w.h.p.

$$L(f_W) \leq \hat{L}(f_W) \cdot \frac{\sum_{i=1}^l \frac{Q_i}{s_i^2} \prod_{i=1}^l r_i^2}{n}$$

- The bound grows exponentially with depth, and is tight in the worst-case

Uniform convergence bounds on Lipschitz-continuous neural networks

Neural network setup

- f_W : A multi-layer neural net with weight matrices $W = [W_1; W_2; \dots; W_l]$
- D : An unknown distribution over feature space times label space $X \times Y$

Norm bounds for neural nets

- Hypothesis space with bounded ℓ_2 and Frobenius norms:

$$H = \left\{ \begin{array}{l} \sum_{i=1}^n k W_1 k_2 \quad s_1; k W_2 k_2 \quad s_2; \dots; k W_l k_2 \quad s_l; \\ k W_1 k_F \quad s_1 r_1; k W_2 k_F \quad s_2 r_2; \dots; k W_l k_F \quad s_l r_l \end{array} \right\} \quad (1)$$

- For any l -layer deep net f_W whose weight matrices W belong to H , with Lipschitz-continuous activation functions at every layer, the following holds w.h.p.

$$L(f_W) \leq \hat{L}(f_W) \cdot \frac{\sum_{i=1}^l \frac{Q_i}{s_i^2} + \frac{P_i}{r_i^2}}{n}$$

- The bound grows exponentially with depth, and is tight in the worst-case

Central thesis: Measure generalization through the Hessian of loss surfaces

Derivation of a Hessian trace measure on the generalization gap

- Main ideas

- Hessian-based regularization for neural net training

Generalization in graph neural networks

- Setup and results

- Numerical results

Problem setup

An l -layer feedforward neural network with weight matrices $W_1; W_2; \dots; W_l$:

$$f_W(x) = \sigma_l \circ \dots \circ \sigma_3 \circ W_3 \circ \sigma_2 \circ W_2 \circ \sigma_1(W_1 x)$$

- Supervised learning: Train weights from random initialization
- Transfer learning: Adapt weights from pretrained foundations models (e.g., fine-tuning)

Problem setup

An l -layer feedforward neural network with weight matrices $W_1; W_2; \dots; W_l$:

$$f_W(x) = \sigma_l \left(\sum_{j=1}^{l-1} W_{j+1} \sigma_j(W_j x) \right)$$

- Supervised learning: Train weights from random initialization
- Transfer learning: Adapt weights from pretrained foundations models (e.g., fine-tuning)

Problem setup

An l -layer feedforward neural network with weight matrices $W_1; W_2; \dots; W_l$:

$$f_W(x) = \sigma_l \left(\sum_{i=1}^l W_i x \right)$$

- Supervised learning: Train weights from random initialization
- Transfer learning: Adapt weights from pretrained foundations models (e.g., fine-tuning)

Distance from initialization affects generalization

The distance between the initialization and the model can affect generalization:
This corresponds to the spectral radius of the hypothesis space [NK19]

Figure 1: Results from Li and Zhang [LZ21], by fine-tuning ResNet-18 on seven image classification tasks

Figure 2: Similar results have appeared in prior work on transfer learning in deep neural networks (folklore)

Our approach based on PAC-Bayes bounds

PAC-Bayes bounds

- $\hat{P} = N(W^{(0)}; \sigma^2 Id)$: prior distribution centered at initialization weights
- $\hat{Q} = N(W^{(\tau)}; \sigma^2 Id)$: posterior distribution centered at trained weights
- With probability at least $1 - \delta$ for any $\epsilon \in (0, 1)$ (see, e.g., McAllester [McA13])

$$\mathbb{E}_{\hat{Q}} [L(f_w)] \leq \mathbb{E}_{\hat{Q}} \left[\hat{L}(f_w) \right] + \frac{\text{KL}(\hat{Q} \parallel \hat{P}) + \log(4n/\delta)}{n} \quad (2)$$

We will use a generalized result, which subsumes the above result by setting \hat{P} appropriately. For any $\epsilon \in (0, 1)$, with probability at least $1 - \delta$ [Cat07]

$$\mathbb{E}_{\hat{Q}} [L(f_w)] \leq \frac{1}{\epsilon} \mathbb{E}_{\hat{Q}} \left[\hat{L}(f_w) \right] + \frac{\text{KL}(\hat{Q} \parallel \hat{P}) + \log(4n/\delta)}{2(1-\epsilon)n} \quad (3)$$

Our approach based on PAC-Bayes bounds

PAC-Bayes bounds

- $\hat{P} = N(W^{(0)}; \sigma^2 I_d)$: prior distribution centered at initialization weights
- $\hat{Q} = N(W^{(\tau)}; \sigma^2 I_d)$: posterior distribution centered at trained weights
- With probability at least $1 - \delta$ for any $\epsilon \in (0, 1)$ (see, e.g., McAllester [McA13])

$$\mathbb{E}_{\hat{Q}} [L(f_w)] \leq \mathbb{E}_{\hat{Q}} [L(f_w)] + \frac{\text{KL}(\hat{Q} \parallel \hat{P}) + \log(4n/\delta)}{n} \quad (2)$$

We will use a generalized result, which subsumes the above result by setting \hat{P} appropriately. For any $\epsilon \in (0, 1)$, with probability at least $1 - \delta$ [Cat07]

$$\mathbb{E}_{\hat{Q}} [L(f_w)] \leq \frac{1}{2} \mathbb{E}_{\hat{Q}} [L(f_w)] + \frac{\text{KL}(\hat{Q} \parallel \hat{P}) + \log(4n/\delta)}{2(1-\epsilon)n} \quad (3)$$

Our approach based on PAC-Bayes bounds

PAC-Bayes bounds

- $P = N(W^{(0)}; \sigma^2 I_d)$: prior distribution centered at initialization weights
- $Q = N(W^{(\tau)}; \sigma^2 I_d)$: posterior distribution centered at trained weights
- With probability at least $1 - \delta$ for any $\eta \in (0; 1)$ (see, e.g., McAllester [McA13])

$$\mathbb{E}_Q [L(f_w)] \leq \mathbb{E}_Q [\hat{L}(f_w)] + \frac{\text{KL}(Q||P) + \log(4n/\delta)}{n} \quad (2)$$

We will use a generalized result, which subsumes the above result by setting η appropriately. For any $\eta \in (0; 1)$, with probability at least $1 - \delta$ [Cat07]

$$\mathbb{E}_Q [L(f_w)] \leq \frac{1}{\eta} \mathbb{E}_Q [\hat{L}(f_w)] + \frac{\text{KL}(Q||P) + \log(4/\delta)}{2(1-\eta)n} \quad (3)$$

Our approach based on PAC-Bayes bounds

PAC-Bayes bounds

- $P = N(W^{(0)}; \sigma^2 I_d)$: prior distribution centered at initialization weights
- $Q = N(W^{(T)}; \sigma^2 I_d)$: posterior distribution centered at trained weights
- With probability at least $1 - \delta$ for any $\beta \in (0; 1)$ (see, e.g., McAllester [McA13])

$$\mathbb{E}_{w \sim Q} [L(f_w)] \leq \mathbb{E}_{w \sim Q} [\hat{L}(f_w)] + \frac{\text{KL}(Q \| P) + \log(4n/\delta)}{n} \quad (2)$$

We will use a generalized result, which subsumes the above result by setting appropriately. For any $\beta \in (0; 1)$, with probability at least $1 - \delta$ [Cat07]

$$\mathbb{E}_{w \sim Q} [L(f_w)] \leq \frac{1}{\beta} \mathbb{E}_{w \sim Q} [\hat{L}(f_w)] + \frac{\text{KL}(Q \| P) + \log(4/\delta)}{2(1-\beta)n} \quad (3)$$

Our finding: the Hessian of loss surface can also affect generalization

Claim 1: Use Taylor's expansion on the perturbed loss

Measure model stability after adding perturbations to weight parameters. Let

$\ell_Q(f_W) = E_{W \sim Q} [\ell(f_W)] = E_{U \sim N(0; \sigma^2 I_D)} [\ell(f_{W+U})]$. We have

$$\ell_Q(f_W(x); y) = \ell(f_W(x); y) + \frac{1}{2} \text{Tr} \left[\sigma^2 \mathbb{H}^2 \ell(f_W(x); y) \right] + C_1 \sigma^3 \quad (4)$$

Our finding: the Hessian of loss surface can also affect generalization

Claim 1: Use Taylor's expansion on the perturbed loss

Measure model stability after adding perturbations to weight parameters. Let

$\ell_Q(f_W) = \mathbb{E}_{W \sim Q} [\ell(f_W)] = \mathbb{E}_{U \sim N(0; \sigma^2 \text{Id})} [\ell(f_{W+U})]$. We have

$$\ell_Q(f_W(x); y) \approx \ell(f_W(x); y) + \frac{1}{2} \text{Tr} \left[\sigma^2 \text{H}^2 \ell(f_W(x); y) \right] + \mathcal{O}(\sigma^3) \quad (4)$$

Figure 3: Illustration of the Hessian approximation in equation (4). We report the results at the last epoch. σ : standard deviation of Gaussian noise injected into weight matrices.

Our finding: the Hessian of loss surface can also affect generalization

Claim 1: Use Taylor's expansion on the perturbed loss

Measure model stability after adding perturbations to weight parameters. Let

$\hat{Q}(f_W) = E_{W \sim Q} [\hat{Q}(f_W)] = E_{U \sim N(0; \sigma^2 I_D)} [\hat{Q}(f_{W+U})]$. We have

$$\hat{Q}(f_W(x); y) - \hat{Q}(f_W(x); y) = \frac{1}{2} \text{Tr} \left[\text{Hess}_{f_W(x); y} [\hat{Q}(f_W(x); y)] \right] \leq C_1 \sigma^2 \quad (4)$$

Claim 2: Uniform convergence of the Hessian operator

Because the Hessian operator is Lipschitz-continuous, the trace of the Hessian satisfies the uniform convergence within the hypothesis space

A sketch for deriving the Hessian trace measure

From equation (3),

$$L_Q(W) = \frac{1}{2} \hat{L}_Q(W) + \frac{C(\text{KL}(Q_{jj}P) + \log(\frac{1}{\epsilon}))}{2(1-\epsilon)n} \quad (5)$$

By Claim 1 (applied to both the expected loss L_Q and the empirical loss \hat{L}_Q),

$$L_Q(W) = L(W) + \frac{\epsilon}{2} \mathbb{E}_{(x,y) \sim D} \text{Tr} \left(r^{-2} \cdot \text{Hess}_{f_W(x); y} \right) + O(\epsilon^3) \quad (6)$$

$$\hat{L}_Q(W) = \hat{L}(W) + \frac{\epsilon}{2n} \sum_{i=1}^n \text{Tr} \left(r^{-2} \cdot \text{Hess}_{f_W(x_i); y_i} \right) + O(\epsilon^3) \quad (7)$$

By Claim 2, the difference between the second terms of equations (6), (7) is of order $O(n^{-\frac{1}{2}})$

By choosing ϵ^2 and r carefully, we get [ZLJ24]

$$L(W) = (1 + \epsilon) \hat{L}(W) + (1 + \epsilon) \frac{C}{n} r^2 + O(n^{-\frac{3}{2}} \log(\frac{1}{\epsilon})) ; \quad (8)$$

where ϵ is the trace norm of the hypothesis space, and r is the radius of the hypothesis space

A sketch for deriving the Hessian trace measure

From equation (3),

$$L_Q(W) = \hat{L}_Q(W) + \frac{C(\text{KL}(Q_{jj}P) + \log(\frac{1}{n}))}{2(1-\gamma)n} \quad (5)$$

By Claim 1 (applied to both the expected loss L_Q and the empirical loss \hat{L}_Q),

$$L_Q(W) = L(W) + \frac{\gamma}{2} \mathbb{E}_{(x;y) \sim D} \text{Tr} \left(r^{2\gamma} (f_W(x); y) \right) + O(n^{-3}) \quad (6)$$

$$\hat{L}_Q(W) = \hat{L}(W) + \frac{\gamma}{2n} \sum_{i=1}^n \text{Tr} \left(r^{2\gamma} (f_W(x_i); y_i) \right) + O(n^{-3}) \quad (7)$$

By Claim 2, the difference between the second terms of equations (6), (7) is of order $O(n^{-\frac{1}{2}})$

By choosing γ and r carefully, we get [ZLJ24]

$$L(W) = (1 + \gamma) \hat{L}(W) + (1 + \gamma) \frac{C r^{2\gamma}}{n} + O(n^{-\frac{3}{4}} \log(\frac{1}{n})) ; \quad (8)$$

where γ is the trace norm of the hypothesis space, and r is the radius of the hypothesis space

A sketch for deriving the Hessian trace measure

From equation (3),

$$L_Q(W) = \hat{L}_Q(W) + \frac{C(\text{KL}(Q_{jj}P) + \log(\frac{1}{\epsilon}))}{2(1-\epsilon)n} \quad (5)$$

By Claim 1 (applied to both the expected loss L_Q and the empirical loss \hat{L}_Q),

$$L_Q(W) = L(W) + \frac{\epsilon^2}{2} \mathbb{E}_{(x,y) \sim D} \text{Tr} \left(r^{2\epsilon} (f_W(x); y) \right) + O(\epsilon^3) \quad (6)$$

$$\hat{L}_Q(W) = \hat{L}(W) + \frac{\epsilon^2}{2n} \sum_{i=1}^n \text{Tr} \left(r^{2\epsilon} (f_W(x_i); y_i) \right) + O(\epsilon^3) \quad (7)$$

By Claim 2, the difference between the second terms of equations (6), (7) is of order $O(n^{-\frac{1}{2}})$

By choosing ϵ^2 and ϵ carefully, we get [ZLJ24]

$$L(W) = (1 + \epsilon) \hat{L}(W) + (1 + \epsilon) \frac{C \cdot r^2}{n} + O(n^{-\frac{3}{4}} \log(\frac{1}{\epsilon})) ; \quad (8)$$

where ϵ is the trace norm of the hypothesis space, and r is the radius of the hypothesis space

A sketch for deriving the Hessian trace measure

From equation (3),

$$L_Q(W) = \frac{1}{2} \hat{L}_Q(W) + \frac{C(\text{KL}(Q_{jj}P) + \log(\frac{1}{\epsilon}))}{2(1-\epsilon)n} \quad (5)$$

By Claim 1 (applied to both the expected loss L_Q and the empirical loss \hat{L}_Q),

$$L_Q(W) = L(W) + \frac{\epsilon}{2} \mathbb{E}_{(x;y) \sim D} \text{Tr} \left(r^{-2} \text{Hess}_{(x;y)}(f_W(x); y) \right) + O(\epsilon^3) \quad (6)$$

$$\hat{L}_Q(W) = \hat{L}(W) + \frac{\epsilon}{2n} \sum_{i=1}^n \text{Tr} \left(r^{-2} \text{Hess}_{(x_i; y_i)}(f_W(x_i); y_i) \right) + O(\epsilon^3) \quad (7)$$

By Claim 2, the difference between the second terms of equations (6), (7) is of order $O(n^{-\frac{1}{2}})$

By choosing ϵ^2 and r carefully, we get [ZLJ24]

$$L(W) = (1 + \epsilon) \hat{L}(W) + (1 + \epsilon) \frac{C r^2}{n} + O(n^{-\frac{3}{4}} \log(\frac{1}{\epsilon})) ; \quad (8)$$

where ϵ is the trace norm of the hypothesis space, and r is the radius of the hypothesis space

Numerical results

The Hessian trace measure correlates well with empirical measurements of the generalization gap. Compare across seven re-tuning methods: SGD, early stopping, weight decay, label smoothing, mixup, distance-based regularization, sharpness-aware minimization

(a) ResNet-50

(b) ResNet-50

(c) BERT-Base

Figure 3: The Hessian measures accurately correlate with empirical generalization errors for seven re-tuning methods

A very related study pursuing a compression-based approach [LFK+22]

Noise injection induces a regularization of the Hessian

An algorithmic implication : Instead of minimizing the loss \hat{L} , we could minimize $\hat{Q} = E_{U \sim N(0, \sigma^2 I_d)} [\hat{L}(f_{W+U})]$ instead. By equation (4), this regularizes trace of the Hessian of the loss surface

Experimental results: Comparison between SGD, noise injection (directly add noise before computing gradient, i.e., WP-SGD), and a two-point noise injection algorithm (NSO)

Noise injection induces a regularization of the Hessian

An algorithmic implication : Instead of minimizing the loss \mathcal{L} , we could minimize $\mathcal{L}_Q = \mathbb{E}_{U \sim \mathcal{N}(0, \sigma^2 \text{Id})} [\mathcal{L}(f_{W+U})]$ instead. By equation (4), this regularizes trace of the Hessian of the loss surface

Experimental results: Comparison between SGD, noise injection (directly add noise before computing gradient, i.e., WP-SGD), and a two-point noise injection algorithm (NSO)

Figure 3: Fine-tuning ResNet-34 and BERT-Base, respectively, on an image and a text classification dataset. Similar results for more recent architectures (multi-modal, CLIP), chain-of-thought fine-tuning (LM, transformer) [ZLJ24]

Derivation of a Hessian trace measure on the generalization gap

Main ideas

Hessian-based regularization for neural net training

Generalization in graph neural networks

Setup and results

Numerical results

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction**: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l - 1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma \left(\sum_{i \in \mathcal{N}(v)} X U^{(t)} + \sum_{i \in \mathcal{N}(v)} P_{G^{(t-1)}}(H^{(t-1)}) W^{(t)} \right); \text{ for } t = 1; 2; \dots; l - 1 \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector v for all $v \in V$: let X be a feature matrix
- Graph-level prediction: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l - 1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma \left(\sum_{i \in \mathcal{N}(v)} X U^{(t)} + P_{G^{(t-1)}}(H^{(t-1)}) W^{(t)} \right); \text{ for } t = 1; 2; \dots; l - 1 \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l-1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma \left(\sum_{i \in \mathcal{N}(v)} X U^{(t)} + P_{G^{(t-1)}}(H^{(t-1)}) W^{(t)} \right); \text{ for } t = 1; 2; \dots; l-1 \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction**: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l-1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma \left(\sum_{i \in \mathcal{N}(v)} X U^{(t)} + P_{G^{(t-1)}}(H^{(t-1)}) W^{(t)} \right); \text{ for } t = 1; 2; \dots; l-1 \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Message-passing neural networks

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction**: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l-1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma_t(XU^{(t)} + \sum_{i \in \mathcal{N}(v)} P_G(v, i)(H^{(t-1)})W^{(t)}); \text{ for } t = 1; 2; \dots; l-1 \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction**: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l - 1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma_t(XU^{(t)} + \sigma_t(P_G \sigma_t(H^{(t-1)}))W^{(t)}); \text{ for } t = 1; 2; \dots; l \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

Graph neural networks are neural networks designed for working with graphs

- Graph $G = (V; E)$: V is a set of vertices and E is a set of edges
- Every node has a feature vector x_v for all $v \in V$: let X be a feature matrix
- Graph-level prediction**: a prediction label $y \in Y$ for each graph instance G

Message passing neural networks (MPNNs) with a pre-defined graph diffusion matrix/mechanism P_G (on each graph G)

- Let $H^{(0)} = X$
- For the first $l - 1$ layers, recursively compute node embedding

$$H^{(t)} = \sigma_t(XU^{(t)} + \sigma_t(P_G \sigma_t(H^{(t-1)}))W^{(t)}); \text{ for } t = 1; 2; \dots; l \quad (9)$$

- For the last layer l , aggregate the embedding of all nodes

$$H^{(l)} = \frac{1}{n} \mathbf{1}_n^T H^{(l-1)} W^{(l)}; \quad (10)$$

How does the generalization bound of GNNs scale with graph structures?

Known results apply to the following GNN architectures

- ^ Graph convolutional networks (GCN) [KW17]
- ^ GraphSAGE (SAmple and aggreGatE) [HYL17]
- ^ Graph isomorphism networks (GIN) [XHL+19]

Table 1: Illustration of the dependence on graph structure of existing results. d : max degree, l : number of GNN layers (depth), A : adjacency matrix, D : degree-diagonal matrix of A

	GCN	MPNN	GIN	GraphSAGE
Garg et al. [GJJ20]	d^{l-1}	d^{l-1}	-	-
Liao et al. [LUZ21]	$d^{\frac{l-1}{2}}$	d^{l-1}	-	-
Ours [JLS+23]	1	$kA k_2^{l-1}$	$\prod_{i=1}^{l-1} \frac{kA k_2^i}{1-1}$	$D^{-1} A^{\frac{l-1}{2}}$

The dependence on graph structure, scaled with maximum degree, for the spectral norm of P_G , which is provably $\leq d$

How does the generalization bound of GNNs scale with graph structures?

Known results apply to the following GNN architectures

- ^ Graph convolutional networks (GCN) [KW17]
- ^ GraphSAGE (SAmple and aggreGatE) [HYL17]
- ^ Graph isomorphism networks (GIN) [XHL+19]

Example on a one-layer linear GNN

One-layer linear GNN with average pooling of node embeddings $\{v_j\}$

$$f(X; G) = \frac{1}{n} \mathbf{1}_n^\top P_G X W^{(1)}$$

Thus, the Euclidean norm of $f(X; G)$ is less than

$$\begin{aligned} \|f(X; G)\| &= \left\| \frac{1}{n} \mathbf{1}_n^\top P_G X W^{(1)} \right\| \\ &\leq \frac{1}{n} \mathbf{1}_n^\top P_G X W^{(1)} := C \end{aligned}$$

Provided that the loss function $\ell(\cdot; y)$ is Lipschitz-continuous, given N samples, we know from standard results that (see, e.g., Zhang [Zha23])

$$L(f) - \hat{L}(f) \leq \frac{C}{N}$$

Example on a one-layer linear GNN

One-layer linear GNN with average pooling of node embeddings $\{v_j\}$

$$f(X; G) = \frac{1}{n} \mathbf{1}_n^\top P_G X W^{(1)}$$

Thus, the Euclidean norm of $f(X; G)$ is less than

$$\begin{aligned} \|f(X; G)\| &= \left\| \frac{1}{n} \mathbf{1}_n^\top P_G X W^{(1)} \right\| \\ &\leq \frac{1}{n} \mathbf{1}_n^\top \left(\|P_G\|_2 \|X\|_2 \|W^{(1)}\|_2 \right) := C \end{aligned}$$

Provided that the loss function $\ell(\cdot; y)$ is Lipschitz-continuous, given N samples, we know from standard results that (see, e.g., Zhang [Zha23])

$$L(f) - \hat{L}(f) \leq \frac{r}{\sqrt{N}}$$

Example on a one-layer linear GNN

One-layer linear GNN with average pooling of node embeddings $\mathbf{g} = \frac{1}{n} \sum_j \mathbf{v}_j$

$$f(\mathbf{X}; \mathbf{G}) = \frac{1}{n} \mathbf{1}_n^\top \mathbf{P}_G \mathbf{X} \mathbf{W}^{(1)}$$

Thus, the Euclidean norm of $f(\mathbf{X}; \mathbf{G})$ is less than

$$\begin{aligned} \|f(\mathbf{X}; \mathbf{G})\| &= \left\| \frac{1}{n} \mathbf{1}_n^\top \mathbf{P}_G \mathbf{X} \mathbf{W}^{(1)} \right\| \\ &\leq \frac{1}{n} \|\mathbf{1}_n\|_2 \|\mathbf{P}_G\|_2 \|\mathbf{X}\|_2 \|\mathbf{W}^{(1)}\|_2 := C \end{aligned}$$

Provided that the loss function $\ell(\cdot; \mathbf{y})$ is Lipschitz-continuous, given N samples, we know from standard results that (see, e.g., Zhang [Zha23])

$$L(f) - \hat{L}(f) \leq \frac{r}{N}$$

Generalizing the simple example to MPNNs

Generalization bounds for MPNNs (informal) [JLS+23]

Suppose activations and loss function are all twice-differentiable, Lipschitz-continuous, first-order and second-order derivatives are both Lipschitz-continuous. d_i : number of neurons at layer i , for $i = 1; 2; \dots; l$

With probability at least $1 - \epsilon$ over N samples, for any $\epsilon > 0$, and any $\epsilon > 0$, any GNN f with weights in H (recall equation (1)) satisfies:

$$L(f) \leq (1 + \epsilon) \hat{L}(f) + \frac{\sum_{i=1}^l d_i \max_{(x;G;y) \in D} \|k_X k_2^2 k_{P_G} k_2^{2(l-1)}\| r_i^2 \sum_{j=1}^Q s_j^2}{N} \quad (11)$$

Example: for GCN, $P_G = D^{-1/2} A D^{-1/2}$, hence $\|k_{P_G} k_2\| \leq 1$

Open problem: remove the dependence on d_i to get **size-independent sample complexity** for GNNs? Known for feedforward NN [GRS18]

Generalizing the simple example to MPNNs

Generalization bounds for MPNNs (informal) [JLS+23]

Suppose activations and loss function are all twice-differentiable, Lipschitz-continuous, first-order and second-order derivatives are both Lipschitz-continuous. d_i : number of neurons at layer i , for $i = 1; 2; \dots; l$

With probability at least $1 - \epsilon$ over N samples, for any $\epsilon > 0$, and any $\epsilon > 0$, any GNN f with weights in H (recall equation (1)) satisfies:

$$L(f) \leq (1 + \epsilon) \hat{L}(f) + \frac{\sum_{i=1}^l d_i \max_{(x; G; y) \in D} \|k_X k_2^2 k_{P_G} k_2^{2(l-1)}\|_2^2 \sum_{j=1}^Q r_i^2 s_j^2}{N} \quad (11)$$

Example: for GCN, $P_G = D^{-1/2} A D^{-1/2}$, hence $\|k_{P_G} k_2\|_2 = 1$

Open problem: remove the dependence on d_i to get size-independent sample complexity for GNNs? Known for feedforward NN [GRS18]

Generalizing the simple example to MPNNs

Generalization bounds for MPNNs (informal) [JLS+23]

Suppose activations and loss function are all twice-differentiable, Lipschitz-continuous, first-order and second-order derivatives are both Lipschitz-continuous. d_i : number of neurons at layer i , for $i = 1; 2; \dots; l$

With probability at least $1 - \epsilon$ over N samples, for any $\epsilon > 0$, and any $\epsilon > 0$, any GNN f with weights in H (recall equation (1)) satisfies:

$$L(f) \leq (1 + \epsilon) \hat{L}(f) + \frac{\sum_{i=1}^l d_i \max_{(x;G;y) \in D} \|k_X k_2^2 k_{P_G} k_2^{2(l-1)}\|_2^2 \sum_{j=1}^Q r_i^2 s_j^2}{N} \quad (11)$$

Example: for GCN, $P_G = D^{-1/2} A D^{-1/2}$, hence $\|k_{P_G} k_2\|_2 = 1$

Open problem: remove the dependence on d_i to get **size-independent sample complexity** for GNNs? Known for feedforward NN [GRS18]

(a) Two-layer GCN

(b) Four-layer GCN

(c) Six-layer GCN

(d) Two-layer MPNN

(e) Four-layer MPNN

(f) Six-layer MPNN

Figure 5: Comparing our result and prior results Garg et al. [GJJ20] and Liao et al. [LUZ21] on three graph classification tasks conducted on GCNs and MPNNs

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Conclusion

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Conclusion

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Conclusion

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Conclusion

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Conclusion

The use of Hessian to study neural networks seems to have been abandoned by the community due to their high computational costs

In this work, we propose to use Hessian of the loss surface as a measure of the generalization gap of neural networks and language models

Some partial progress

- ^ Can explain a variety of phenomena observed with neural network training
- ^ Provides a new approach to prove generalization bounds for GNNs

Next steps

- ^ Better understand the structure of Hessian of loss surfaces in large models
- ^ Known sample complexity for GNN only works for graph-level prediction; what about the sample complexity for learning node-level prediction tasks

Thank you for listening!

Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. (2019). "Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks". In: International Conference on Machine Learning PMLR, pp. 322-332 (6{13}).

Arora, S., Ge, R., Neyshabur, B., and Zhang, Y. (2018). "Stronger generalization bounds for deep nets via a compression approach". In: International Conference on Machine Learning PMLR, pp. 254-263 (6{13}).

Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. (2017). "Spectrally-normalized margin bounds for neural networks". In: Advances in neural information processing systems 30 (6{13}).

Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. (2020). "Benign overfitting in linear regression". In: Proceedings of the National Academy of Sciences 117.48, pp. 30063-30070 (6{13}).

Catoni, O. (2007). "PAC-Bayesian supervised classification: the thermodynamics of statistical learning". In: arXiv preprint arXiv:0712.0248 (23{26}).

Garg, V., Jegelka, S., and Jaakkola, T. (2020). "Generalization and representational limits of graph neural networks". In: ICML (46, 54).

Golowich, N., Rakhlin, A., and Shamir, O. (2018). "Size-independent sample complexity of neural networks". In: Conference On Learning Theory PMLR, pp. 297{299 (51{53).

Hamilton, W., Ying, Z., and Leskovec, J. (2017). "Inductive representation learning on large graphs". In: NeurIPS (46, 47).




Ju, H., Li, D., Sharma, A., and Zhang, H. R. (2023). "Generalization in graph neural networks: Improved pac-bayesian bounds on graph diffusion". In: International Conference on Artificial Intelligence and Statistics. PMLR, pp. 6314{6341 (46, 51{53).

Ju, H., Li, D., and Zhang, H. R. (2022). "Robust re-tuning of deep neural networks with hessian-based generalization guarantees". In: International Conference on Machine Learning PMLR, pp. 10431{10461 (6{13).

Kipf, T. N. and Welling, M. (2017). "Semi-supervised classification with graph convolutional networks". In: ICLR (46, 47).

Li, D. and Zhang, H. (2021). "Improved regularization and robustness for re-tuning in neural networks". In: Advances in Neural Information Processing Systems34, pp. 27249{27262 (22).

-  Li, Y., Ma, T., and Zhang, H. (2018). **“Algorithmic Regularization in Over-parameterized Matrix Sensing and Neural Networks with Quadratic Activations”**. In: *Conference On Learning Theory* (6–13).
-  Liao, R., Urtasun, R., and Zemel, R. (2021). **“A PAC-Bayesian Approach to Generalization Bounds for Graph Neural Networks”**. In: *ICLR* (46, 54).
-  Lot , S., Finzi, M., Kapoor, S., Potapczynski, A., Goldblum, M., and Wilson, A. G. (2022). **“PAC-Bayes compression bounds so tight that they can explain generalization”**. In: *Advances in Neural Information Processing Systems 35*, pp. 31459–31473 (35).
-  McAllester, D. (2013). **“A PAC-Bayesian tutorial with a dropout bound”**. In: *arXiv preprint arXiv:1307.2118* (23–26).
-  Nagarajan, V. and Kolter, J. Z. (2019). **“Uniform convergence may be unable to explain generalization in deep learning”**. In: *Advances in Neural Information Processing Systems 32* (22).
-  Shannon, C. E. (1948). **“A mathematical theory of communication”**. In: *The Bell system technical journal* 27.3, pp. 379–423 (5).
-  Xu, K., Hu, W., Leskovec, J., and Jegelka, S. (2019). **“How powerful are graph neural networks?”** In: *ICLR* (46, 47).

-  Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. (2021). **“Understanding deep learning (still) requires rethinking generalization”**. In: *Communications of the ACM* 64.3, pp. 107–115 (3).
-  Zhang, H. R., Li, D., and Ju, H. (2024). **“Noise Stability Optimization for Finding Flat Minima: A Hessian-based Regularization Approach”**. In: *Transactions on Machine Learning Research* (30–33, 37).
-  Zhang, T. (2023). **Mathematical analysis of machine learning algorithms**. Cambridge University Press (6–13, 48–50).