Generalization Theory for Deep Learning



Basic version: finite hypothesis class

Finite hypothesis class: with probability $1 - \delta$ over the choice of a training set of size n, for a bounded loss ℓ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x_i), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x_i, y_i) \sim D} \left[\ell(f(x_i), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x_i, y_i) \sim D} \left[\ell(f(x_i), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x_i, y_i) \sim D} \left[\ell(f(x_i), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}\right)$$

$$\lim_{f \in \mathcal{F}} \ell(f(x_i), y_i) - \mathbb{E}_{(x_i, y_i) \sim D} \left[\ell(f(x_i), y_i) \right] \right| = O\left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/$$

VC-Dimension

Motivation: Do we need to consider **every** classifier in \mathcal{F} ? Intuitively, **pattern of classifications** on the training set should suffice. (Two predictors that predict identically on the training set should generalize similarly).

Let $\mathcal{F} = \{f : \mathbb{R}^d \to \{+1, -1\}\}$ be a class of binary classifiers.

The growth function $\Pi_{\mathscr{F}}: \mathbb{N} \to \mathbb{F}$ is defined as:

$$\Pi_{\mathcal{F}}(m) = \max_{(x_1, x_2, \dots, x_m)} \left| \left\{ (f(x_1), f(x_2), \dots, f(x_m)) \mid f \in \mathcal{F} \right\} \right|.$$

The VC dimension of \mathcal{F} is defined as:

$$VCdim(\mathcal{F}) = \max\{m : \Pi_{\mathcal{F}}(m) = 2^m\}.$$

VC-dimension Generalization bound

Theorem (Vapnik-Chervonenkis): with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

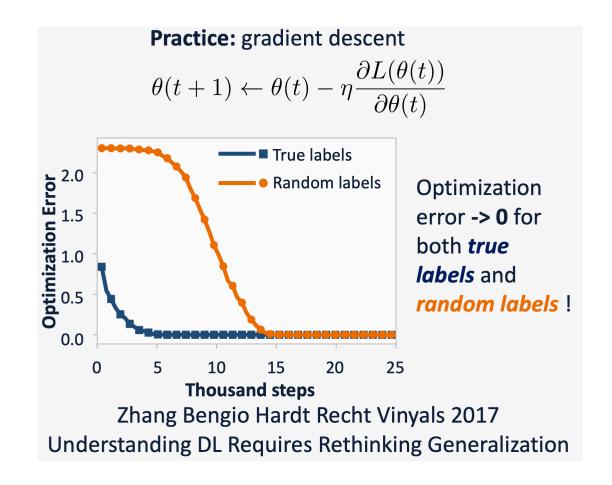
$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\mathsf{VCdim}(\mathcal{F}) \log n + \log 1/\delta}{n}}\right)$$

Examples:

- Linear functions: VC-dim = O(dimension)
- Neural network: VC-dimension of fully-connected net with width W and H layers is $\Theta(WH)$ (Bartlett et al., '17).

Problems with VC-dimension bound $\mu < \# \mu$

- 1. In over-parameterized regime, bound >> 1.
- 2. Cannot explain the random noise phenomenon:
 - Neural networks that fit random labels and that fit true labels have the same VC-dimension.



PAC Bayesian Generalization Bounds

Setup: Let P be a prior over function in class \mathcal{F} , let Q be the posterior (after algorithm's training).

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{KL(Q \mid \mid P) + \log 1/\delta}{n}}\right)$$

$$\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{O}(\mathcal{O}) : \mathcal{O}(\mathcal{$$

Rademacher Complexity

Intuition: how well can a classifier class fit random noise?

(Empirical) Rademacher complexity: For a training set $S = \{x_1, x_2, ..., x_n\}$, and a class \mathcal{F} , denote:

$$\hat{R}_n(S) = \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(x_i) .$$

where $\sigma_i \sim \text{Unif}\{+1, -1\}$ (Rademacher R.V.).

(Population) Rademacher complexity:

$$R_n = \mathbb{E}_S \left[\hat{R}_n(s) \right].$$

Rademacher Complexity Generalization Bound

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss \mathcal{E} , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\frac{\hat{R}_n}{n} + \sqrt{\frac{\log 1/\delta}{n}}\right)$$

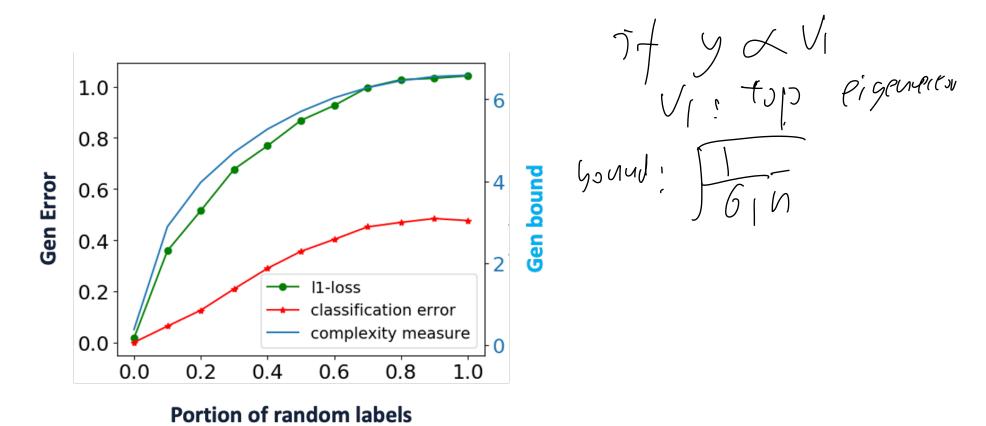
and

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\frac{R_n}{n} + \sqrt{\frac{\log 1/\delta}{n}}\right)$$

Kernel generalization bound

7: 2 h

Use Rademacher complexity theory, we can obtain a generalization bound $O(\sqrt{y^{\mathsf{T}}(H^*)^{-1}y/n})$ where $y \in \mathbb{R}^n$ are n labels, and $H^* \in \mathbb{R}^{n \times n}$ is the kernel (e.g., NTK) matrix.

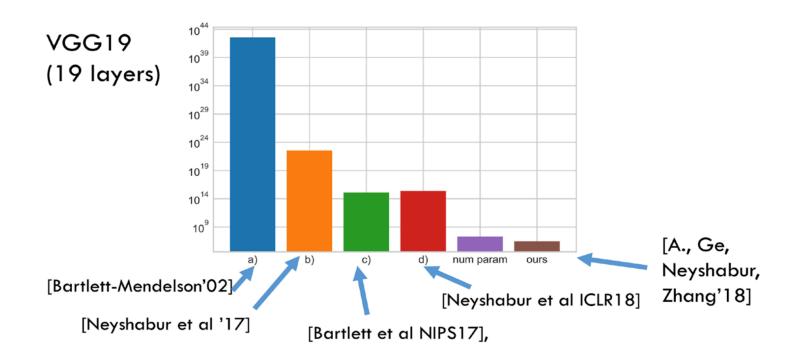


Norm-based Rademacher complexity bound

Theorem: If the activation function is σ is ρ -Lipschitz. Let $\mathscr{F} = \{x \mapsto W_{H+1}\sigma(W_h\sigma(\cdots\sigma(W_1x)\cdots),\|W_h^T\|_{1,\infty} \leq B \,\forall h \in [H]\}$ then $R_n(\mathscr{S}) \leq \|X^T\|_{2,\infty} (2\rho B)^{H+1} \sqrt{2\ln d}$ where $X = [x_1,\ldots,x_n] \in \mathbb{R}^{d\times n}$ is the input data matrix.

Comments on generalization bounds

- When plugged in real values, the bounds are rarely non-trivial (i.e., smaller than 1)
- "Fantastic Generalization Measures and Where to Find them" by Jiang et al. '19: large-scale investigation of the correlation of extant generalization measures with true generalization.



Comments on generalization bounds

- Uniform convergence may be unable to explain generalization of deep learning [Nagarajan and Kolter, '19]
 - Uniform convergence: a bound for all $f \in \mathcal{F}$
 - Exists example that 1) can generalize, 2) uniform convergence fails.
- Rates:
 - Most bounds: $1/\sqrt{n}$.
 - Local Rademacher complexity: 1/n.

 For approximation and optimization, neural network has no advantage over kernel. Why NN gives better performance: generalization.

• [Allen-Zhu and Li '20] Construct a class of functions ${\mathcal F}$ such that y = f(x) for some $f \in \mathcal{F}$:
• no kernel is sample-efficient;

Doly (d) Sample/

- Exists a neural network that is sample-efficient.

Separation between NN and kernel Thin 2 a day of functions (S, (2) 2) and a distuibution over 2 d:MS-t, i) I kernel method, if it satisfies that H (E C) GNPM 45- ((M)) if $E_{M} \left[\left(C(x) - Cf, \Phi(x) \right) \right] \stackrel{\text{def}}{=} \frac{1}{2}$ then you need n > 2Ti) = Simple procedure (ast kernel) that you can output correct c, using u=d samples ANN can Simulate this procedure

part i) C is a basis for
$$\{f: \{f, I\} \rightarrow R\}$$

w. v.t. W

(s, Cs' $E_{X-M} \left(C_{S}(X) \cdot C_{S}(X) \right) = \{0\} \text{ if } S \neq S \}$
 $= \} \forall f, \exists \{ m \}, f = \exists a \cdot C_{S}$

Goal: $E_{X-M} \left(C_{S}(X) \cdot C_{S}(X) \right) = \{1\} \text{ if } S = S \}$

some $f = \underbrace{\exists a \cdot e_{X}(X)}_{S} \text{ if } S = \underbrace{\exists a \cdot e$

Exam
$$[(s^*(x) - cf, \phi(x))^2]$$

= $[(s^*(x) - 2f, \phi(x))^2]$

Notort Nus;

$$\Lambda : 2^{\alpha} \times 9$$

$$\Lambda : 3^{\alpha} : = 3^{\alpha} / 5$$

$$\Lambda : = 0 \times 2^{\alpha}$$

Separation between NN and kernel 52 = diag (SL) + 52', 52': off-diagond $|(\mathcal{N}')|_{\overline{f}}^2 = \sum_{g} e_{ig}e_{g}^2(\mathcal{N}') \leq \frac{2^d}{q}$ =) SL has at most = 29 eigenvalue 7/3 consider subspace with eigenvalue 2/3 I N C. Hair subspace UX C this subspace [[SLX [[] = [[diag (s) x+ six [] 2 >/ (Idiag (D) x 1/2 - (D2 'x 1/2 Vauk(SL) 7 = 2d 7 = 1/1/2 - 3/1/2 - 3/1/2 > 0