



# Theoretical Analysis of Weak-to-Strong Generalization

Neurips 2024

Hunter Lang, David Sontag, Aravindan Vijayaraghavan

# Problem Setup

## Notation

- $\mathbf{x} \sim \mathcal{D}$ , assume input space  $\mathcal{X}$  is discrete.
- Ground-truth function  $y : \mathcal{X} \rightarrow \mathcal{Y} = \{1, \dots, k\}$
- Pseudo-labeler (Teacher)  $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y} \cup \{\emptyset\} \rightarrow \text{Label}$  is either in  $\mathcal{Y}$  or does not exist.
- Covered set  $S = \{\mathbf{x} | \tilde{y}(\mathbf{x}) \neq \emptyset\}$ : subset of  $\mathcal{X}$  that has a pseudo-label.
- Uncovered set  $T = \{\mathbf{x} | \tilde{y}(\mathbf{x}) = \emptyset\} = \mathcal{X} \setminus S$
- Partitioning based on the ground-truth:  $\mathcal{X}_i = \{\mathbf{x} | y(\mathbf{x}) = i\}$ ,  $S_i = S \cap \mathcal{X}_i$ ,  $T_i = T \cap \mathcal{X}_i \rightarrow \{\mathcal{X}_i\}$ ,  $\{S_i\}$ , and  $\{T_i\}$  are partitions of  $\mathcal{X}$ ,  $S$ , and  $T$  respectively.
- Partitioning  $S_i$ :  $S_i^{good} : \{\mathbf{x} \in S_i | \tilde{y}(\mathbf{x}) = y(\mathbf{x})\}$ ,  $S_i^{bad} : \{\mathbf{x} \in S_i | \tilde{y}(\mathbf{x}) \neq y(\mathbf{x})\}$
- Error rate of pseudo-labeler:  $\alpha_i = \mathbb{P}(S_i^{bad} | S_i)$ , assume  $\alpha_i \in (0, 1/2)$  for all  $i$ .
- Strong model hypothesis class:  $\mathcal{F}$

## Problem Setup.

- Training on the covered set targeted by the weak label:  $f^* = \arg \min_{\mathcal{F}} \text{err}(f, \tilde{y}|S)$ .
- Goal: Obtain upper bounds on the error:  $\text{err}(f^*, y|\mathcal{X})$ .

# Illustrative Example: Sentiment Classifier

## Settings

- $\mathcal{X}$ : Text documents,  $\mathcal{Y} = \{-1, +1\}$ .
- Pseudo-label

$$\tilde{y} = \begin{cases} +1, & \text{if 'incredible' } \in x, \\ -1, & \text{if 'horrible' } \in x, \\ \emptyset, & \text{otherwise.} \end{cases}$$

Assume 'incredible' and 'horrible' never co-occur.

- $\alpha_{-1} = \alpha_{+1} = \alpha$ .

## Existing Error bounds

- Proposition 3.1 (Bound from [1]):  $\tilde{f} = \mathbb{P}(S) \cdot 4\alpha(1 - \alpha) + \mathbb{P}(T)$   
→ No pseudo-label correction in  $\alpha < 3/4$ .
- Proposition 3.2 (Bound from [2]):  $\mathbb{P}(S)$  should be larger than  $2/3$ .  
→ Cannot explain weak-to-strong generalization in the low-coverage regime.  
⇒ Existing studies fail to explain weak-to-strong generalization.

# Neighborhood and $\eta$ -Robustness

## Neighborhood $\mathcal{N}$

- General definition: Function that maps each point  $\mathbf{x}$  to a set of points  $\mathcal{N}(\mathbf{x}) \in \mathcal{X}$ .
- Symmetric assumption:  $\mathbf{x} \in \mathcal{N}(\mathbf{x}') \Leftrightarrow \mathbf{x}' \in \mathcal{N}(\mathbf{x})$ .
- Examples:  $\mathcal{N}(\mathbf{x}) = \{\mathbf{x}' : \|\varphi(\mathbf{x}) - \varphi(\mathbf{x}')\| \leq r\}$  for some rep. function  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$ .

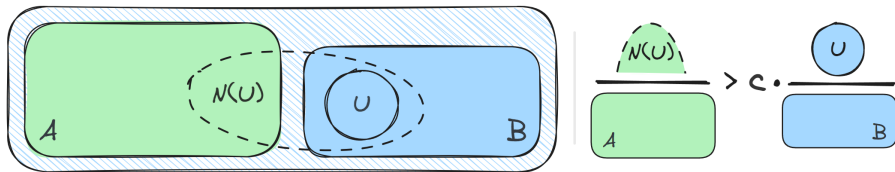
## $\eta$ -Robustness

- $r(f, \mathbf{x}) = \mathbb{P}(f(\mathbf{x}') \neq f(\mathbf{x}) | \mathbf{x}' \in \mathcal{N}(\mathbf{x}))$
- $\eta$ -robust at  $\mathbf{x}$ :  $r(f, \mathbf{x}) \leq \eta$ .
- $R_\eta(f) = \{\mathbf{x} : r(f, \mathbf{x}) \leq \eta\}$ : set of  $\eta$ -robust points for  $f$ .
- Average-case robustness: Classifier gives the same labels for most of their neighbors.

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{x}' \sim \mathcal{D} | \mathcal{N}(\mathbf{x})} [f(\mathbf{x}) \neq f(\mathbf{x}')] \leq \gamma$$

is  $\eta$ -robust in probability at least  $1 - \gamma/\eta$ .

# Expansion



**Expansion:** For fixed sets  $A, B \in \mathcal{X}$ ,  $\mathbb{P}_x$  satisfies  $(c, q)$ -expansion on  $(A, B)$  if

$$\forall U \in B \text{ with } \mathbb{P}(U|B) > q, \quad \mathbb{P}(\mathcal{N}(U)|A) > c\mathbb{P}(U|B).$$

**Expansion of a set collection** (Relaxed version): Suppose  $\mathcal{M}$  is a collection of subsets of  $B$ .  $\mathcal{M}$  is  $(c, q)$ -expansion on  $(A, B)$  if

$$\forall U \in \mathcal{M} \text{ with } \mathbb{P}(U|B) > q, \quad \mathbb{P}(\mathcal{N}(U)|A) > c\mathbb{P}(U|B).$$

# Adversarially Robust Models ( $\eta = 0$ )

## Expansion

- "Bad" points (incorrect pseudolabels) have many "good" neighbors.
- "Uncovered" points (no pseudolabels) have many "good" neighbors.

Idea: Student model with "robust" on the neighborhoods

### Theorem 4.1 (Pseudo-label correction, informal.)

The true error of  $f$  on covered set  $S_i$  satisfies:

$$\text{err}(f, y|S_i) \leq \text{err}(f, \tilde{y}|S_i) + \alpha_i \left(1 - \frac{3}{2}c\right).$$

Trivial bounds:  $\text{err}(f, y|S_i) \leq \text{err}(f, \tilde{y}|S_i) + \alpha_i \rightarrow$  tighter than trivial bounds.

### Theorem 4.2 (Coverage Expansion, informal.)

The true error of  $f$  on uncovered set  $T_i$  satisfies:

$$\text{err}(f, y|T_i) \leq \text{err}(\bar{R}(f)|T_i) + \max \left( q, \frac{\text{err}(f, \tilde{y}|S_i) - c\alpha_i}{c(1 - 2\alpha_i)} \right).$$

# References

- [1] D. Fu, M. Chen, F. Sala, S. Hooper, K. Fatahalian, and C. Ré, “Fast and three-rious: Speeding up weak supervision with triplet methods,” in *International conference on machine learning*, pp. 3280–3291, PMLR, 2020.
- [2] C. Wei, K. Shen, Y. Chen, and T. Ma, “Theoretical analysis of self-training with deep networks on unlabeled data,” *arXiv preprint arXiv:2010.03622*, 2020.